

Abstract We use the H -matrix technology to compute the approximate square root of a covariance matrix in linear cost. This allows us to generate normal and log-normal random fields on general point sets with optimal cost. We derive rigorous error estimates which show convergence of the method. Our approach requires only mild assumptions on the covariance function and on the point set. Therefore, it might be also a nice alternative to the circulant embedding approach which applies only to regular grids and stationary covariance functions.

Fast random field generation with H -matrices

Michael Feischl · Frances Kuo · Ian H. Sloan

Received: date / Accepted: date

1 Introduction

Generating samples of random fields is a common bottleneck in simulation and modeling of real life phenomena as, e.g., structural vibrations [6], groundwater flow [8], and composite material behavior [1]. A standard approach is to truncate the Karhunen-Loève expansion of the random field. This can, particularly for rough fields with short correlation length, be very expensive, as many summands of the expansion have to be evaluated to compute a decent approximation. Often, it suffices to evaluate the random field only on some particular (quadrature) nodes. If the random field $a(\mathbf{x}, \omega)$ is normally or log-normally distributed, it is well-known that the evaluation at the quadrature nodes can be done by factorizing the corresponding covariance matrix \mathbf{C} of the normal random field as $\mathbf{C} = \mathbf{A}\mathbf{A}^T$ for some matrix \mathbf{A} . Since each evaluation requires a matrix-vector multiplication with \mathbf{A} , the direct approach requires $\mathcal{O}(N^2)$ operations, while the factorization itself requires $\mathcal{O}(N^3)$ operations and thus is prohibitively expensive. An efficient method first proposed in [4,3] is *circulant embedding*, which employs fast FFT techniques to realize the factorization and the matrix-vector multiplication in $\mathcal{O}(N \log(N))$ operations. This approach, however, works solely for stationary covariance functions $\varrho(\mathbf{x}, \mathbf{y}) = \rho(|\mathbf{x} - \mathbf{y}|)$ and regular grids of quadrature nodes. Since non-stationary covariance functions are of great interest for the modeling of natural structures (e.g., porous rock, wood, ...), and since finite element methods often use irregular grids, we propose a new method which removes both restrictions.

The idea is to approximate the covariance matrix \mathbf{C} by an H^2 -matrix, as described in, e.g., [2], and to use an iterative method to approximate $\mathbf{A}\mathbf{z}$ for any $\mathbf{z} \in \mathbb{R}^N$. This is feasible since matrix-vector multiplication with H^2 -matrices can be done in $\mathcal{O}(N)$ operations. The only assumption on the covariance function of the random field is that it is asymptotically smooth. We propose two iterative algorithms, each with individual advantages for smooth or rough random fields. This algorithms might also be of interest for the approximation of random fields with covariance kernels of random solutions of certain stochastic operator equations, as considered in [5].

The idea to use H -matrices for random field approximation has already been used indirectly in [10], where the authors efficiently compute Eigenfunctions of the covariance operator by use of H -matrix techniques.

1.1 Notation

Throughout the text, $\alpha \lesssim \beta$ denotes $\alpha \leq C\beta$ for some generic constant $C > 0$ and $\alpha \simeq \beta$ means $\alpha \lesssim \beta$ and $\beta \lesssim \alpha$. The notation $|\cdot|$ has several unambiguous meanings: for vectors, it denotes the euclidean norm, while for sets, $|\cdot|$ is the natural measure, which is the Lebesgue measure (volume, area) for continuous sets and the counting measure (cardinality) for finite sets. The notation $\|\cdot\|_2$ is used for the

M. Feischl, F Kuo, I.H. Sloan
School of Mathematics and Statistics, UNSW Sydney, NSW 2052
Tel.: +61-2-93857076
E-mail: m.feischl@unsw.edu.au

spectral matrix norm and $|\mathbf{z}|_p := (\sum_{j=1}^N |z_j|^p)^{1/p}$ for all $\mathbf{z} \in \mathbb{R}^N$ denotes the ℓ_p -norm. By \mathcal{P}^k we denote the set of polynomials of maximal degree k . For brevity, we write $|\cdot| := |\cdot|_2$. We denote the maximal and minimal eigenvalues of a positive definite and symmetric matrix $\mathbf{M} \in \mathbb{R}^{N \times N}$ by

$$\lambda_{\max}(\mathbf{M}) := \sup_{\mathbf{z} \in \mathbb{R}^N \setminus \{0\}} \frac{|\mathbf{M}\mathbf{z}|}{|\mathbf{z}|} \quad \text{and} \quad \lambda_{\min}(\mathbf{M}) := \inf_{\mathbf{z} \in \mathbb{R}^N \setminus \{0\}} \frac{(\mathbf{M}\mathbf{z})^T \mathbf{z}}{|\mathbf{z}|^2}.$$

We denote the k -th component of a vector $\mathbf{v} \in \mathbb{R}^N$ by v_k , whereas sequences of vectors are denoted by $\mathbf{v}^1, \mathbf{v}^2, \dots$.

2 Model Problem

Let $(\Omega, \Sigma, \mathbb{P})$ be a probability space and let $D \subseteq \mathbb{R}^d$, $d \in \mathbb{N}$ be a Lipschitz domain. We consider a random field which is normal or log-normal,

$$\mathcal{Z}(\mathbf{x}, \omega) \quad \text{or} \quad \exp(\mathcal{Z}(\mathbf{x}, \omega)) \quad \text{for all } \omega \in \Omega, \mathbf{x} \in D$$

for some zero-mean Gaussian random field $\mathcal{Z}(\cdot, \cdot)$ (note that the assumption on the mean is purely for brevity of presentation). The covariance function $\varrho: D \times D \rightarrow \mathbb{R}$ of $\mathcal{Z}(\cdot, \cdot)$ is assumed asymptotically smooth: that is, $\varrho \in C^\infty(\{(\mathbf{x}, \mathbf{y}) \in D \times D : \mathbf{x} \neq \mathbf{y}\})$ and there exist constants $c_1, c_2 > 0$ such that

$$|\partial_{\mathbf{x}}^\alpha \partial_{\mathbf{y}}^\beta \varrho(\mathbf{x}, \mathbf{y})| \leq c_1 (c_2 |\mathbf{x} - \mathbf{y}|)^{-|\alpha|_1 - |\beta|_1} |\alpha + \beta|_1! \quad \text{for all } \mathbf{x} \neq \mathbf{y} \in D, \quad (1)$$

for all multi-indices $\alpha, \beta \in \mathbb{N}_0^d$ with $|\alpha|_1 + |\beta|_1 \geq 1$. (The expert reader will notice that the original definition of asymptotically smooth includes a singularity order. As our covariance functions are always finite in value, we do not consider this.) The goal of this work is to derive an efficient method which evaluates the random field at certain (quadrature) points $\mathcal{N} \subseteq D$, where $\mathcal{N} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ is a finite set, i.e., we aim to approximate

$$\left(\mathcal{Z}(\mathbf{x}, \omega) \right)_{\mathbf{x} \in \mathcal{N}} \in \mathbb{R}^N \quad \text{or} \quad \left(\exp(\mathcal{Z}(\mathbf{x}, \omega)) \right)_{\mathbf{x} \in \mathcal{N}} \in \mathbb{R}^N$$

for given $\omega \in \Omega$.

2.1 Examples of valid covariance functions

The condition above includes the important class of isotropic stationary covariance functions of Matérn form, e.g.,

$$\varrho(\mathbf{x}, \mathbf{y}) = \sigma^2 \frac{2^{1-\mu}}{\Gamma(\mu)} \left(\sqrt{2\mu} \frac{|\mathbf{x} - \mathbf{y}|_p}{\lambda} \right)^\mu K_\mu \left(\sqrt{2\mu} \frac{|\mathbf{x} - \mathbf{y}|_p}{\lambda} \right), \quad (2)$$

where $\Gamma(\cdot)$ is the gamma function, K_μ is the modified Bessel function of second kind, and $\lambda, \sigma > 0$, $\mu \in (0, \infty]$, $p \in \mathbb{N}$ are parameters. For $\mu = 1/2$, the above function takes the form

$$\varrho(\mathbf{x}, \mathbf{y}) = \sigma^2 \exp \left(- \frac{|\mathbf{x} - \mathbf{y}|_p}{\lambda} \right)$$

and the limit case $\mu = \infty$ satisfies

$$\varrho(\mathbf{x}, \mathbf{y}) = \sigma^2 \exp \left(- \frac{|\mathbf{x} - \mathbf{y}|_p^2}{2\lambda^2} \right).$$

Also much more general non-stationary, non-isotropic covariance functions, e.g.,

$$\varrho(\mathbf{x}, \mathbf{y}) := \sigma^2 \frac{\det(\boldsymbol{\Sigma}_{\mathbf{x}})^{1/4} \det(\boldsymbol{\Sigma}_{\mathbf{y}})^{1/4}}{\sqrt{2} \det(\boldsymbol{\Sigma}_{\mathbf{x}} + \boldsymbol{\Sigma}_{\mathbf{y}})^{1/2}} \exp \left(- (\mathbf{x} - \mathbf{y})^T \frac{(\boldsymbol{\Sigma}_{\mathbf{x}} + \boldsymbol{\Sigma}_{\mathbf{y}})^{-1}}{2} (\mathbf{x} - \mathbf{y}) \right). \quad (3)$$

satisfy the assumptions. Here, $\boldsymbol{\Sigma}(\cdot): D \rightarrow \mathbb{R}^{d \times d}$ is a smooth mapping into the symmetric positive definite matrices and $\sigma > 0$ is a parameter. This covariance function was first suggested in [11] to model spatially dependent anisotropies in a material.

Lemma 1 *The covariance functions from (2) satisfy (1). Assume the mapping $\mathbf{x} \mapsto \boldsymbol{\Sigma}_{\mathbf{x}}$ satisfies (for any matrix norm $\|\cdot\|$)*

$$\sup_{\alpha \in \mathbb{N}^d} \sup_{\mathbf{x} \in D} \|\partial_{\mathbf{x}}^{\alpha} \boldsymbol{\Sigma}_{\mathbf{x}}\| < \infty. \quad (4)$$

Then, the covariance function from (3) is asymptotically smooth (1).

We postpone the proof of the lemma to Appendix A.

3 Sampling the random field

By definition, $\mathcal{Z}(\mathbf{x}, \cdot)$, $\mathbf{x} \in \mathcal{N}$ is a Gaussian random field with covariance matrix $\mathbf{C} \in \mathbb{R}^{N \times N}$, $N = |\mathcal{N}|$, and $\mathbf{C}_{ij} = \varrho(\mathbf{x}_i, \mathbf{x}_j)$, where we write $\mathcal{N} := \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$. The main goal of this section is to establish a new way to efficiently approximate $\mathbf{C}^{1/2} \mathbf{z}$ for given $\mathbf{z} \in \mathbb{R}^N$. Roughly, the strategy is to approximate \mathbf{C} by an H^2 -matrix and to benefit from the fast matrix-vector multiplication provided by it. This allows us to efficiently approximate $\mathbf{A} \mathbf{z}$ (without actually factorizing the matrix \mathbf{C}).

3.1 H^2 -matrix approximation of the covariance matrix

Given the finite set of evaluation points $\mathcal{N} := \{\mathbf{x}_1, \dots, \mathbf{x}_N\} \subset D$, we approximate the covariance matrix $\mathbf{C} \in \mathbb{R}^{N \times N}$, $\mathbf{C}_{ij} := \varrho(\mathbf{x}_i, \mathbf{x}_j)$ by an H^2 -matrix \mathbf{C}_p via interpolation of order $p \in \mathbb{N}$.

In the following, we recall the definition of H^2 -matrices and the approximation process as laid out in, e.g., [2]. The rough idea is to partition the index set of the covariance matrix into *far-field* blocks, which can be approximated efficiently by interpolation of the covariance function, and *near-field* blocks, which are stored exactly.

3.1.1 Block partitioning

For each subset $X \subseteq \mathcal{N}$, we denote by $B_X \subseteq \mathbb{R}^d$, the smallest axis-parallel box such that $X \subseteq B_X$. We build a binary tree of clusters in the following way. Let $X_{\text{root}} := \mathcal{N} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ denote the root of the tree which has level zero $\text{level}(X_{\text{root}}) = 0$ by definition. For each node of the tree X with $|X| > C_{\text{leaf}}$ for some cut-off constant $C_{\text{leaf}} \in \mathbb{N}$ (usually $C_{\text{leaf}} \approx 20$), we define two sons of X as follows: Split B_X in half along its longest edge into $B_0 \cup B_1 = B_X$. Define $\text{sons}(X) := \{X_0, X_1\}$ with $X_0 := X \cap B_0$ and $X_1 := X \setminus X_0$ and set $\text{level}(X_i) = \text{level}(X) + 1$ for $i = 0, 1$. For a node X with $|X| \leq C_{\text{leaf}}$, we define $\text{sons}(X) := \emptyset$. This procedure generates a binary tree denoted by \mathbb{T}_{cl} (where cl stands for *cluster*) and guarantees that its leaves satisfy $|X| \leq C_{\text{leaf}}$.

For a parameter $\eta > 0$, we consider the admissibility condition for axis parallel boxes $B, B' \subseteq \mathbb{R}^d$

$$\max\{\text{diam}(B), \text{diam}(B')\} \leq \eta \text{dist}(B, B'), \quad (5)$$

where the euclidean distance between the bounding boxes is defined by

$$\text{dist}(B, B') := \inf_{\mathbf{x} \in B, \mathbf{y} \in B'} |\mathbf{x} - \mathbf{y}|.$$

The condition (5) will be used to build the block-cluster tree $\mathbb{T} \subseteq \mathbb{T}_{\text{cl}} \times \mathbb{T}_{\text{cl}}$ as follows. The root of \mathbb{T} is $(X_{\text{root}}, X_{\text{root}})$. For each node $(X, Y) \in \mathbb{T}$ of the tree, define $\text{sons}(X, Y)$, the set of sons, as:

$$\begin{cases} \text{if } B_X \text{ and } B_Y \text{ satisfy (5) or if } \text{sons}(X) = \emptyset = \text{sons}(Y) & \text{set } \text{sons}(X, Y) = \emptyset \\ \text{else if } \text{sons}(Y) \neq \emptyset \text{ and } \text{sons}(X) = \emptyset & \text{set } \text{sons}(X, Y) = \{X\} \times \text{sons}(Y) \\ \text{else if } \text{sons}(X) \neq \emptyset \text{ and } \text{sons}(Y) = \emptyset & \text{set } \text{sons}(X, Y) = \text{sons}(X) \times \{Y\} \\ \text{else } \text{sons}(X) \neq \emptyset \text{ and } \text{sons}(Y) \neq \emptyset & \text{set } \text{sons}(X, Y) = \text{sons}(X) \times \text{sons}(Y) \end{cases}$$

We also define the level as $\text{level}(X_{\text{root}}, X_{\text{root}}) = 0$ and $\text{level}(X, Y) = \text{level}(X', Y') + 1$ for $(X, Y) \in \text{sons}(X', Y')$. Further, we define

$$\mathbb{T}_{\text{far}} := \{(X, Y) \in \mathbb{T} : \text{sons}(X, Y) = \emptyset \text{ and } B_X, B_Y \text{ satisfy (5)}\}$$

as well as

$$\mathbb{T}_{\text{near}} := \{(X, Y) \in \mathbb{T} : \text{sons}(X, Y) = \emptyset \text{ and } B_X, B_Y \text{ do not satisfy (5)}\}.$$

Note that by definition of the block-cluster tree \mathbb{T} , the set $\mathbb{T}_{\text{near}} \cup \mathbb{T}_{\text{far}}$ contains all the leaves of \mathbb{T} . Moreover, we see that for each $(X, Y) \in \mathbb{T} \setminus (\mathbb{T}_{\text{near}} \cup \mathbb{T}_{\text{far}})$, there holds

$$X \times Y = \bigcup_{(X', Y') \in \text{sons}(X, Y)} X' \times Y'$$

Therefore, $\mathbb{T}_{\text{near}} \cup \mathbb{T}_{\text{far}}$ is a partition of $\mathcal{N} \times \mathcal{N}$ in the sense that each pair of points $(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{N} \times \mathcal{N}$ for $1 \leq i, j \leq N$ is contained in exactly one $(X, Y) \in \mathbb{T}_{\text{near}} \cup \mathbb{T}_{\text{far}}$.

3.1.2 Interpolation

The blocks $(X, Y) \in \mathbb{T}_{\text{far}}$ satisfy (5) and hence interpolation of the kernel function is highly accurate. This allows us to store the matrix very efficiently. Let $I(X) := \{i \in \mathbb{N} : \mathbf{x}_i \in X\}$ denote the index set of X . The basic idea now is to replace $\mathbf{C}|_{I(X) \times I(Y)}$ by a low-rank approximation $\mathbf{V}^X \mathbf{M}^{XY} (\mathbf{V}^Y)^T$ with $\mathbf{V}^X \in \mathbb{R}^{|X| \times p^d}$, $\mathbf{M}^{XY} \in \mathbb{R}^{p^d \times p^d}$, and $\mathbf{V}^Y \in \mathbb{R}^{|Y| \times p^d}$, where p is the interpolation order. The three matrices are defined by Cebyshev interpolation of the covariance function. To that end, let $\{q_1^X, \dots, q_{p^d}^X\}$ denote transformed, tensorial Cebyshev nodes in B_X with the corresponding Lagrange basis functions $L_1^X, \dots, L_{p^d}^X : B_X \rightarrow \mathbb{R}$. Given $(X, Y) \in \mathbb{T}_{\text{far}}$, we may approximate

$$\varrho(\mathbf{x}, \mathbf{y}) \approx c_p^{XY}(\mathbf{x}, \mathbf{y}) := \sum_{n, m=1}^{p^d} \varrho(q_n^X, q_m^Y) L_n^X(\mathbf{x}) L_m^Y(\mathbf{y}) \quad \text{for all } \mathbf{x} \in X, \mathbf{y} \in Y.$$

For $i, j \in \{1, \dots, N\}$ and $n, m \in \{1, \dots, p^d\}$, this leads to

$$\mathbf{V}_{in}^X := L_n^X(\mathbf{x}_i), \quad \mathbf{V}_{jm}^Y := L_m^Y(\mathbf{x}_j), \quad \text{and} \quad \mathbf{M}_{nm}^{XY} := \varrho(q_n^X, q_m^Y)$$

and hence

$$\mathbf{C}|_{I(X) \times I(Y)} \approx \mathbf{V}^X \mathbf{M}^{XY} (\mathbf{V}^Y)^T.$$

The admissibility condition (5) guarantees that the approximation error converges to zero exponentially in p , as we prove in Proposition 1 below. Further note that the Cebyshev interpolation described above is exact on polynomials of degree p . Thus, for $X \in \mathbb{T}_{\text{cl}}$ and $\mathbf{x}_i \in X' \in \text{sons}(X)$, there holds with the transfer matrices $\mathbf{T}^{X'X} := (L_n^X(q_m^{X'}))_{mn} \in \mathbb{R}^{p^d \times p^d}$

$$\mathbf{V}_{in}^X := L_n^X(\mathbf{x}_i) = \sum_{m=1}^{p^d} L_n^X(q_m^{X'}) L_m^{X'}(\mathbf{x}_i) = \sum_{m=1}^{p^d} L_n^X(q_m^{X'}) \mathbf{V}_{im}^{X'} = (\mathbf{V}^{X'} \mathbf{T}^{X'X})_{in}.$$

Thus, it suffices to store \mathbf{V}^X only for the leaves of \mathbb{T}_{cl} together with the transfer matrices $\mathbf{T}^{X'X}$. This enables very efficient storage and arithmetics for H^2 matrices.

The capabilities of H^2 -matrices which we employ in this work are summarized below in Proposition 1. To that end, we assume that the points \mathcal{N} are approximately uniformly distributed, in the following sense.

Assumption 1 (quasi-uniform distribution) *We say that \mathcal{N} is quasi-uniformly distributed if for all $\mathbf{x} \in \mathcal{N}$, there exists a compact subset $Q_{\mathbf{x}} \subseteq D$ with $\mathbf{x} \in Q_{\mathbf{x}}$ such that $\bigcup_{\mathbf{x} \in \mathcal{N}} Q_{\mathbf{x}} = D$, $Q_{\mathbf{x}} \cap Q_{\mathbf{x}'}$ has measure zero for all $\mathbf{x} \neq \mathbf{x}' \in \mathcal{N}$, and*

$$\begin{aligned} C_u^{-1} |D| &\leq |Q_{\mathbf{x}}| N \leq C_u |D|, \\ C_u^{-1} |Q_{\mathbf{x}}|^{1/d} &\leq \text{diam}(Q_{\mathbf{x}}) \leq C_u |Q_{\mathbf{x}}|^{1/d} \end{aligned} \tag{6}$$

for some constant $C_u > 0$ independent of N .

Proposition 1 Suppose we have a covariance matrix $\mathbf{C} \in \mathbb{R}^{N \times N}$ and an asymptotically smooth kernel $g(\cdot, \cdot)$ and recall Assumption 1 on approximate uniform distribution of \mathcal{N} . Then, there exists a constant $C_H > 0$ such that, for all $p \in \mathbb{N}_0$, the H^2 -matrix $\mathbf{C}_p \in \mathbb{R}^{N \times N}$ constructed as above satisfies

$$\|\mathbf{C} - \mathbf{C}_p\|_2 \leq \|\mathbf{C} - \mathbf{C}_p\|_F := \left(\sum_{i,j=1}^N |\mathbf{C} - \mathbf{C}_p|_{ij}^2 \right)^{1/2} \leq C_H N (\log(p) + 1)^{2d-1} \left(\frac{\eta}{4c_2} \right)^p. \quad (7)$$

(The constant c_2 is defined in (1).) The H^2 -matrix \mathbf{C}_p is symmetric and can be stored using less than $C_H p^{2d} N$ memory units. Moreover, given any vector $\mathbf{x} \in \mathbb{R}^N$, it is possible to compute $\mathbf{C}_p \mathbf{x} \in \mathbb{R}^N$ in less than $C_H p^{2d} N$ arithmetic operations. The constant C_H depends only on C_{leaf} and d . The matrix \mathbf{C}_p is positive definite if p is sufficiently large such that

$$C_H N (\log(p) + 1)^{2d-1} \left(\frac{\eta}{4c_2} \right)^p < \lambda_{\min}(\mathbf{C}). \quad (8)$$

We postpone the proof of the lemma to Appendix B.

3.2 Computing the square-root (Method 1)

Since \mathbf{C} is positive definite in our case, a standard method is to compute the Cholesky factorization $\mathbf{L}\mathbf{L}^T = \mathbf{C}$. This can be done using H^2 -matrices in linear cost as analyzed in [9]. However, to the author's best knowledge, there is no complete error analysis available, and due to the complicated structure of the algorithm, the worst-case error estimate may be overly pessimistic. Therefore, we propose an iterative algorithm based on a variant of the Lanczos iteration. Note that polynomial or rational approximations of the square root (as pursued in, e.g., [15]) are doomed to fail since smooth random fields result in very badly conditioned covariance matrices \mathbf{C} (see also the numerical experiments below). This implies that a polynomial approximation of the square root over the spectrum of \mathbf{C} is very costly, whereas a rational approximation requires the inverse of \mathbf{C} which is hard to compute due to the bad condition number.

The idea behind the algorithm below is as follows. Given a positive definite symmetric matrix $\mathbf{M} \in \mathbb{R}^{N \times N}$ and a vector $\mathbf{z} \in \mathbb{R}^N$, the aim is to compute efficiently an approximation to $\mathbf{M}^{1/2} \mathbf{z}$. For arbitrary $k \leq N$ define the order- k Krylov subspace of \mathbf{M} and \mathbf{z} as

$$\mathcal{K}_k := \text{span}\{\mathbf{z}, \mathbf{M}\mathbf{z}, \mathbf{M}^2\mathbf{z}, \dots, \mathbf{M}^{(k-1)}\mathbf{z}\}. \quad (9)$$

Assuming \mathcal{K}_k is k -dimensional, consider the orthogonal matrix $\mathbf{Q} \in \mathbb{R}^{N \times k}$ whose columns are the orthonormal basis vectors of the Krylov subspace, i.e., $\mathbf{Q}^T \mathbf{Q} = \mathbf{I}_k$ and $\text{range}(\mathbf{Q}) = \mathcal{K}_k$. Now define $\mathbf{U} \in \mathbb{R}^{k \times k}$ by

$$\mathbf{U} := \mathbf{Q}^T \mathbf{M} \mathbf{Q}.$$

If $k = N$ then $\mathbf{Q}\mathbf{Q}^T = \mathbf{I}_N$ and $\mathbf{Q}\mathbf{U}\mathbf{Q}^T = \mathbf{M}$, from which it follows that

$$\mathbf{M}^{1/2} \mathbf{z} = \mathbf{Q} \mathbf{U}^{1/2} \mathbf{Q}^T \mathbf{z}. \quad (10)$$

The algorithm relies on explicit matrix multiplication to construct \mathbf{U} and then a direct factorization of \mathbf{U} , thus for large N it is feasible only when $k \ll N$, in which case (10) does not hold exactly. However, as we show later it may hold to a good enough approximation. The following Lanczos type algorithm builds up progressively the columns of \mathbf{Q} without fully computing \mathcal{K}_k first.

Algorithm 1 *Input:* positive definite symmetric matrix $\mathbf{M} \in \mathbb{R}^{N \times N}$, vector $\mathbf{z} \in \mathbb{R}^N$, and maximal number of iterations $k \in \mathbb{N}$.

1. Compute Krylov subspace: Set $\mathbf{Q}_1 := \mathbf{z}/|\mathbf{z}| \in \mathbb{R}^{N \times 1}$ and $k_0 = k$. For $j = 2, \dots, k$ do:
 - (a) Compute $\tilde{\mathbf{q}} := \mathbf{M}\mathbf{q}^{j-1} \in \mathbb{R}^N$, where \mathbf{q}^{j-1} is the $(j-1)$ -th column of $\mathbf{Q}_{j-1} \in \mathbb{R}^{N \times (j-1)}$.
 - (b) Compute QR-factorization $\mathbf{Q}_j \in \mathbb{R}^{N \times j}$ (with orthonormal columns), $\mathbf{R}_j \in \mathbb{R}^{j \times j}$ (upper triangular) such that $\mathbf{Q}_j \mathbf{R}_j = (\mathbf{Q}_{j-1}, \tilde{\mathbf{q}}) \in \mathbb{R}^{N \times j}$.
 - (c) If $(\mathbf{R}_j)_{jj} = 0$, set $k_0 = j-1$ and goto (2).
2. Compute $\mathbf{U}_{k_0} := \mathbf{Q}_{k_0}^T \mathbf{M} \mathbf{Q}_{k_0} \in \mathbb{R}^{k_0 \times k_0}$.
3. Compute $\mathbf{U}_{k_0}^{1/2}$ directly.

4. Return $\mathbf{y} = \mathbf{Q}_{k_0} \mathbf{U}_{k_0}^{1/2} \mathbf{Q}_{k_0}^T \mathbf{z}$.

Output: Approximation $\mathbf{y} \approx \mathbf{M}^{1/2} \mathbf{z}$ and number of steps k_0 .

Remark 1 Obviously, the orthogonal basis $\mathbf{q}^1, \dots, \mathbf{q}^k$ could also be generated by Gram-Schmidt orthogonalization. However, numerical experiments show that this is not stable with respect to roundoff errors. Moreover, also the classical Lanczos algorithm seems to be prone to rounding errors, especially for ill-conditioned matrices. Therefore, we propose to use the QR -factorization as above.

Remark 2 As proved in Lemma 5 below (and as is easily verified), \mathbf{Q}_j contains the first j columns of \mathbf{Q} , hence the first j columns of \mathbf{Q}_{j+1} coincide with \mathbf{Q}_j . Thus, it suffices to store only the new column \mathbf{q}^j .

Theorem 1 Let $0 < \eta < 4c_2$ and let p be sufficiently large such that \mathbf{C}_p constructed from \mathbf{C} as in Section 3.1 is positive definite (condition (8) is sufficient), and suppose Assumption 1 holds. Given $\mathbf{z} \in \mathbb{R}^N$, call Algorithm 1 with $\mathbf{M} = \mathbf{C}_p$, \mathbf{z} , and a maximal number of iterations $k \in \mathbb{N}$. The output of Algorithm 1 contains the approximation $\mathcal{Z}_{k,p}(\mathbf{z}) := \mathbf{y} \in \mathbb{R}^N$ to $\mathbf{C}^{1/2} \mathbf{z}$ and the step number $k_0 \leq k$.

(i) There holds with Kronecker's delta $\delta_{i,j}$

$$\begin{aligned} & \frac{|\mathbf{C}^{1/2} \mathbf{z} - \mathcal{Z}_{k,p}(\mathbf{z})|}{|\mathbf{z}|} \\ & \leq \delta_{k_0,k} \sqrt{2} \sqrt{\lambda_{\min}(\mathbf{C}_p) + \lambda_{\max}(\mathbf{C}_p)} \left(1 - \frac{2\lambda_{\min}(\mathbf{C}_p)}{\lambda_{\min}(\mathbf{C}_p) + \lambda_{\max}(\mathbf{C}_p)}\right)^{k^{\log_3(2)/2}} \\ & \quad + \frac{2C_H N (\log(p) + 1)^{2d-1} \left(\frac{\eta}{4c_2}\right)^p}{\max\{\lambda_{\min}(\mathbf{C}), \lambda_{\min}(\mathbf{C}_p)\}^{1/2}}. \end{aligned}$$

(ii) Let $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_M > 0$ denote the distinct eigenvalues of \mathbf{C}_p for some $M \leq N$ and assume

$$|\lambda_i - \lambda_j| \leq \lambda_{\max}(\mathbf{C}_p) C_\kappa \kappa^{\min\{i,j\}} \quad \text{for all } 1 \leq i, j \leq M$$

for some $C_\kappa > 0$ and $0 < \kappa < 1$, then

$$\frac{|\mathbf{C}^{1/2} \mathbf{z} - \mathcal{Z}_{k,p}(\mathbf{z})|}{|\mathbf{z}|} \leq C_{\text{alg1}} \delta_{k_0,k} \frac{C_\kappa \lambda_{\max}(\mathbf{C}_p) N^{1/(2k)}}{\sqrt{\lambda_{\min}(\mathbf{C}_p)}} \kappa^{(k+1)/2} + \frac{2C_H N (\log(p) + 1)^{2d-1} \left(\frac{\eta}{4c_2}\right)^p}{\max\{\lambda_{\min}(\mathbf{C}), \lambda_{\min}(\mathbf{C}_p)\}^{1/2}}.$$

The algorithm completes in $\mathcal{O}(k^3 p^{2d} N)$ arithmetic operations and uses less than $\mathcal{O}(kN)$ storage. The constant C_H is defined in Proposition 1 and $C_{\text{alg1}} > 0$ does not depend on \mathbf{C} , p , k , or N .

Remark 3 The theorem covers two regimes of covariance matrices. Whereas case (i) is the classical Lanczos convergence analysis for well-conditioned matrices, case (ii) considers ill-conditioned matrices with rapidly decaying eigenvalues. The numerical examples in Section 4 suggest that the error estimates might be more or less sharp, since Algorithm 1 performs remarkably well for smooth random fields (with rapidly decaying eigenvalues) and very rough random fields (with well-conditioned covariance matrices). Note that $k_0 < k$ (hence $\delta_{k_0,k} = 0$) implies that the condition in the if-clause (1c) is true. This however is an exotic case, meaning that \mathbf{C}_p has some non-trivial invariant subspaces with less than k dimensions. In this situation the algorithm computes $\mathbf{C}_p^{1/2} \mathbf{z}$ exactly and only a H -matrix approximation error remains.

Proof (Proof of Theorem 1) The cost estimate is proved as follows. The Krylov subspace loop of Algorithm 1 completes at most k iterations. In each iteration, we have one H^2 -matrix-vector multiplication which needs $\mathcal{O}(p^{2d} N)$ operations. Moreover, the QR -factorization needs $\mathcal{O}(Nk^2)$ arithmetic operations. After the matrix \mathbf{Q}_k is set up, we have k H^2 -matrix-vector multiplications to compute $\mathbf{M} \mathbf{Q}_k$ and k^2 scalar products to compute \mathbf{U}_{k_0} . In total, this needs $\mathcal{O}(N(k + k^2))$ arithmetic operations. The computation of $\mathbf{U}_{k_0}^{1/2}$ can be done in $\mathcal{O}(k^3)$ operations (see, e.g., [12] for the algorithm and the corresponding analysis). Finally, to compute \mathbf{y} , we have k scalar products, a matrix vector multiplication with a $(k \times k)$ matrix and a matrix-matrix multiplication of $(N \times k)$ and $(k \times k)$ matrices, all of which can be done in $\mathcal{O}(Nk^2)$ arithmetic operations.

To see (i), we employ the triangle inequality

$$\begin{aligned} \frac{|C_p^{1/2}z - Z_{k,p}(z)|}{|z|} &\leq \frac{|C_p^{1/2}z - Z_{k,p}(z)|}{|z|} + \frac{|C_p^{1/2}z - C_p^{1/2}z|}{|z|} \\ &\leq \frac{|C_p^{1/2}z - Z_{k,p}(z)|}{|z|} + \|C_p^{1/2} - C^{1/2}\|_2. \end{aligned} \quad (11)$$

For the first term on the right-hand side, Lemma 7 below proves

$$\frac{|C_p^{1/2}z - Z_{k,p}(z)|}{|z|} \leq \delta_{k_0,k} \sqrt{2} \sqrt{\lambda_{\min}(C_p) + \lambda_{\max}(C_p)} \left(1 - \frac{2\lambda_{\min}(C_p)}{\lambda_{\min}(C_p) + \lambda_{\max}(C_p)}\right)^{k^{\log_3(2)}/2}.$$

As shown in Lemma 2 below, the second term on the right-hand side of (11) is bounded by

$$\|C_p^{1/2} - C^{1/2}\|_2 \leq 2 \max\{\lambda_{\min}(C), \lambda_{\min}(C_p)\}^{-1/2} \|C_p - C\|_2. \quad (12)$$

Hence, (i) follows from Proposition 1. For (ii), we note that the combination of both estimates in Proposition 2 shows for $U_j := Q_j^T M Q_j$

$$\min_{1 \leq j \leq k} \frac{|C_p^{1/2}z - Q_j(U_j^{1/2})Q_j^T z|}{|z|} \leq \delta_{k_0,k} 8C_R^{1/k} \frac{\lambda_{\max}(C_p)N^{1/(2k)}}{\sqrt{\lambda_{\min}(C_p)}} C_\kappa \kappa^{(k+1)/2}.$$

We may eliminate the minimum in the error estimate since Algorithm 1 is essentially (up to roundoff errors) of Lanczos type, and for this algorithm, [7, Example 5.1] shows that the approximation error $|C_p^{1/2}z - Q_j(U_j^{1/2})Q_j^T z|$ decreases monotonically in j . The remainder of the proof then follows as for (i) since $Q_{k_0}(U_{k_0}^{1/2})Q_{k_0}^T z = Z_{k,p}(z)$.

3.3 Computing the square-root (Method 2)

The main drawback of Algorithm 1 is the additional storage requirements due to the necessity to store the matrix Q_k . For this reason, we here follow a different approach, proposing a second algorithm that improves this situation.

The matrix sign function is defined for all square matrices \widetilde{M} with no pure imaginary eigenvalues as

$$\text{sgn}(\widetilde{M}) := \widetilde{M}(\widetilde{M}^2)^{-1/2}.$$

The sign function $\text{sgn}(\widetilde{M})$ can be computed using the Schultz iteration via

$$M_{k+1} = \frac{1}{2}M_k(3I - M_k^2), \quad M_0 = \widetilde{M}. \quad (13)$$

The iterates M_k converge quadratically towards $\text{sgn}(\widetilde{M})$ if $\|I - \widetilde{M}^2\|_2 < 1$ in any matrix norm (see [14, Theorem 5.2]). It is observed in [13], that all matrices $M \in \mathbb{R}^{N \times N}$ with only positive real eigenvalues satisfy

$$\text{sgn} \begin{pmatrix} 0 & M \\ I & 0 \end{pmatrix} = \begin{pmatrix} 0 & M^{1/2} \\ M^{-1/2} & 0 \end{pmatrix},$$

where $I \in \mathbb{R}^{N \times N}$ denotes the identity matrix, which opens the possibility to compute $M^{1/2}$ via the sign function of the matrix. By inserting

$$\widetilde{M} := \begin{pmatrix} 0 & M \\ I & 0 \end{pmatrix}.$$

By inserting this choice of \widetilde{M} into (13), we see that all iterates have the form

$$M_k := \begin{pmatrix} 0 & A_k \\ B_k & 0 \end{pmatrix}.$$

As already observed in [13], this leads to the iteration

$$\mathbf{A}_{k+1} = \frac{1}{2}\mathbf{A}_k(3\mathbf{I} - \mathbf{B}_k\mathbf{A}_k), \quad \mathbf{B}_{k+1} = \frac{1}{2}\mathbf{B}_k(3\mathbf{I} - \mathbf{A}_k\mathbf{B}_k), \quad (14)$$

starting with $\mathbf{A}_0 = \mathbf{M}$ and $\mathbf{B}_0 = \mathbf{I} \in \mathbb{R}^{N \times N}$. The iterates \mathbf{A}_k converge towards $\mathbf{M}^{1/2}$, which is what we aim to compute. The considerations above lead us to the following recursive form of the Schulz algorithm above, which uses only matrix vector multiplication.

Algorithm 2 *Input:* positive definite symmetric matrix $\mathbf{M} \in \mathbb{R}^{N \times N}$, vector $\mathbf{z} \in \mathbb{R}^N$, maximal number of iterations $k \in \mathbb{N}$, temporary storage vectors $\mathbf{z}^j \in \mathbb{R}^N$, $j \in \{1, \dots, k\}$, and scaling factor $0 < s < 2\|\mathbf{C}_p\|_2^{-1}$ (the scaling factor ensures convergence of the algorithm).

Main:

1. Compute $\mathbf{y} = \text{PartA}(s\mathbf{M}, \mathbf{z}, (\mathbf{z}^j)_{j=1}^k, k)$.
2. Return \mathbf{y}/\sqrt{s} .

Output: the approximation $\mathbf{y} \approx \mathbf{M}^{1/2}\mathbf{z}$.

Subroutines:

PartA($\mathbf{M}, \mathbf{z}, (\mathbf{z}^j), k$):

- (i) If $k = 0$, return $\mathbf{M}\mathbf{z}$.
- (ii) Compute $\mathbf{z}^k := \text{PartA}(\mathbf{M}, \mathbf{z}, (\mathbf{z}^j)_{j=1}^{k-1}, k-1)$ and $\mathbf{z}^k := \text{PartB}(\mathbf{M}, \mathbf{z}^k, (\mathbf{z}^j)_{j=1}^{k-1}, k-1)$.
- (iii) Compute $\mathbf{z} := 3\mathbf{z} - \mathbf{z}^k$.
- (iv) Return $\frac{1}{2}\text{PartA}(\mathbf{M}, \mathbf{z}, (\mathbf{z}^j)_{j=1}^{k-1}, k-1)$.

PartB($\mathbf{M}, \mathbf{z}, (\mathbf{z}^j), k$):

- (i) If $k = 0$, return \mathbf{z} .
- (ii) Compute $\mathbf{z}^k := \text{PartB}(\mathbf{M}, \mathbf{z}, (\mathbf{z}^j)_{j=1}^{k-1}, k-1)$ and $\mathbf{z}^k := \text{PartA}(\mathbf{M}, \mathbf{z}^k, (\mathbf{z}^j)_{j=1}^{k-1}, k-1)$.
- (iii) Compute $\mathbf{z} := 3\mathbf{z} - \mathbf{z}^k$.
- (iv) Return $\frac{1}{2}\text{PartB}(\mathbf{M}, \mathbf{z}, (\mathbf{z}^j)_{j=1}^{k-1}, k-1)$.

Theorem 2 Suppose Assumption 1 holds and let $\mathbf{z} \in \mathbb{R}^N$. If $0 < \eta < 4c_2$ and p is sufficiently large such that \mathbf{C}_p constructed from \mathbf{C} as in Section 3.1 is positive definite (condition (8) is sufficient), Algorithm 1 called with $\mathbf{M} = \mathbf{C}_p$ and $0 < s < 2\|\mathbf{C}_p\|_2^{-1}$ computes the approximation $\mathcal{Z}_{k,p}(\mathbf{z}) := \mathbf{y} \in \mathbb{R}^N$ such that

$$\frac{|\mathbf{C}^{1/2}\mathbf{z} - \mathcal{Z}_{k,p}(\mathbf{z})|}{|\mathbf{z}|} \leq s^{-1/2}\kappa^{2^k} + \frac{2C_H N(\log(p) + 1)^{2d-1} \left(\frac{\eta}{4c_2}\right)^p}{\max\{\lambda_{\min}(\mathbf{C}), \lambda_{\min}(\mathbf{C}_p)\}^{1/2}},$$

where $\kappa := \max\{|1 - s\lambda_{\max}(\mathbf{C}_p)|, |1 - s\lambda_{\min}(\mathbf{C}_p)|\} < 1$. The algorithm completes in $\mathcal{O}(3^k p^{2d} N)$ arithmetic operations and uses less than kN extra storage. The constant C_H is defined in Proposition 1.

Remark 4 In contrast to Algorithm 1 which needs $\mathcal{O}(|\log(\varepsilon)|N)$ extra storage (at least in case (ii)), we see that Algorithm 2 requires only $\mathcal{O}(\log|\log(\varepsilon)|N)$ additional storage for an error request of $\varepsilon > 0$.

Proof (Proof of Theorem 2) First, we show that the additional storage vectors are accessed properly without overwriting of necessary information. We prove this by induction on k . For $k = 0$, both **PartA** and **PartB** do not access the extra storage vectors. We confirm that in both parts of the algorithm, the vector \mathbf{z}^j is only used if $k = j$ for the last parameter of **PartA** and **PartB**. During the execution of **PartA**(\cdot, \cdot, \cdot, j) and **PartB**(\cdot, \cdot, \cdot, j), only instances of those functions with $k < j$ are called. By the induction hypothesis, those access the storage vectors properly. This concludes the induction and implies that each instance has the correct value of \mathbf{z}^j .

The computational cost estimate follows directly from Proposition 1 and the fact that each subroutine **PartA** and **PartB** makes at most three function calls. To see the error estimate, we use (11) and note that Algorithm 2 is nothing else then a recursive version of the iteration (14). The scaling $s < 2\|\mathbf{C}_p\|_2^{-1}$ ensures $\kappa < 1$, since $\lambda \in \{\lambda_{\min}(\mathbf{C}_p), \lambda_{\max}(\mathbf{C}_p)\}$ satisfies $1 - s\lambda < 1$ (since $s, \lambda > 0$) as well as $s\lambda - 1 \leq s\|\mathbf{C}_p\|_2 - 1 < 2 - 1 = 1$. Thus, Lemma 9 shows

$$\frac{|\mathbf{C}_p^{1/2}\mathbf{z} - \mathcal{Z}_{k,p}(\mathbf{z})|}{|\mathbf{z}|} \leq s^{-1/2} \left(\max\{|1 - s\lambda_{\max}(\mathbf{C}_p)|, |1 - s\lambda_{\min}(\mathbf{C}_p)|\} \right)^{2^k} = s^{-1/2}\kappa^{2^k}.$$

We conclude the proof with the aid of (12) and Proposition 1.

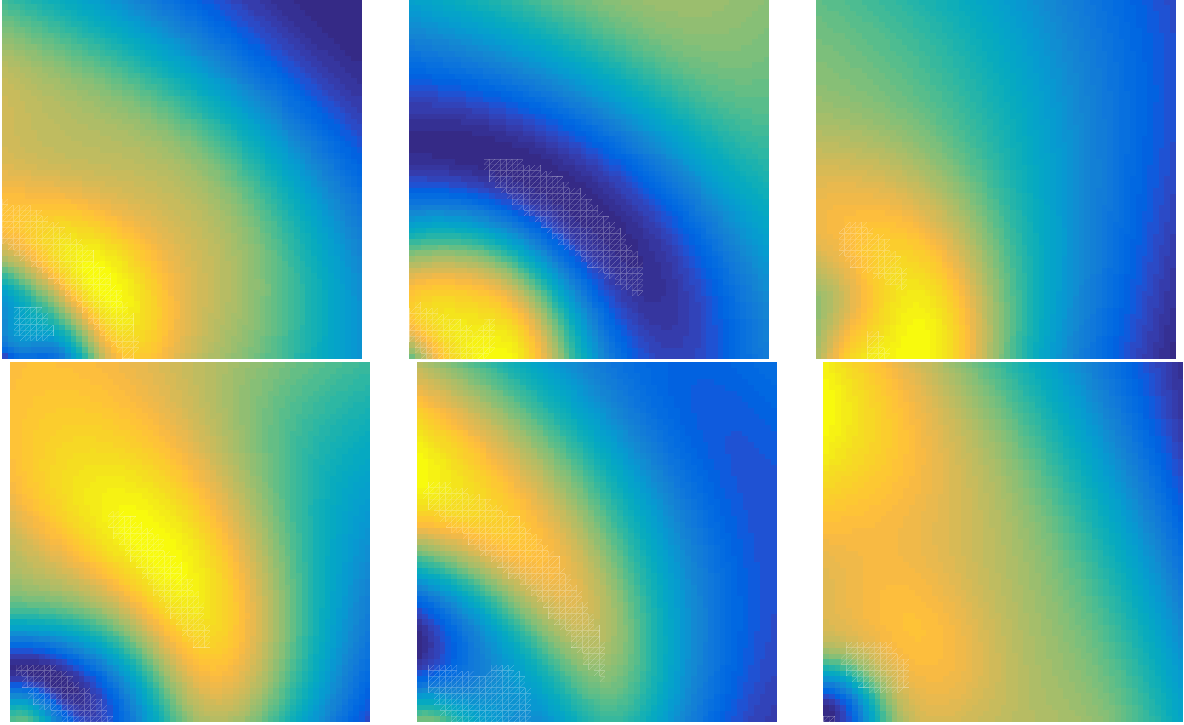


Fig. 1: Samples of \mathcal{Z} with a non-stationary covariance function. We clearly observe the shorter covariance length (more variation) near the bottom left corner.

$m =$	5	6	7	8	9	10
$\lambda = 1$	2.0e+09	6.1e+16	8.6e+17	2.6e+19	1.8e+20	1.4e+20
$\lambda = 10^{-1}$	3.9e+07	5.5e+14	1.8e+17	8.4e+18	4.8e+20	4.6e+20
$\lambda = 10^{-2}$	6.5e+06	2.6e+12	2.7e+17	1.2e+19	3.3e+19	2.8e+20
$\lambda = 10^{-3}$	4.2e+06	9.4e+11	6.1e+17	4.2e+18	2.6e+19	1.1e+20

Table 1: Condition numbers of \mathbf{C} for the covariance function from (2) with \mathcal{N} being a Sobol point set with 2^m points.

4 Numerical experiments

All numerical experiments were computed in Matlab, by use of a Matlab- H^2 -matrix library which can be downloaded under `software.michaelfeischl.net`. The authors are well aware that the Matlab implementation prohibits high-end performance. However, we wanted to demonstrate the feasibility of our algorithms and show the correct convergence rates, for which purpose the Matlab implementation is sufficient.

For the first example, we consider a covariance function of the form (3) with

$$\Sigma_{\mathbf{x}} := |\mathbf{x}|^2 \mathbf{I} \quad \text{and} \quad \Sigma_{\mathbf{y}} := |\mathbf{y}|^2 \mathbf{I}. \quad (15)$$

We use Algorithm 1 to generate six samples on the unit square $D = [0, 1]^2$ of the corresponding normal random field \mathcal{Z} shown in Figure 1. Figure 2–3 show samples of the covariance functions from (2) with different parameters.

To illustrate the challenging nature of handling these covariance matrices, Table 1 shows condition numbers of \mathbf{C} for different problem sizes and the Matérn covariance function (2).

For a performance comparison of Algorithm 1 and Algorithm 12, we consider the covariance function of the form (2) with $p = 2$, $\sigma = 1$, and varying $\mu \in \{1/2, \infty\}$, $\lambda \in \{1, 10^{-1}, 10^{-2}, 10^{-3}\}$. We compute samples of $\mathcal{Z}(\mathbf{x}, \omega)$ on a Sobol pointset with 2^{10} points. The results are plotted in Figure 4 where we see the relative approximation error versus the computation time in seconds. We observe that with respect to

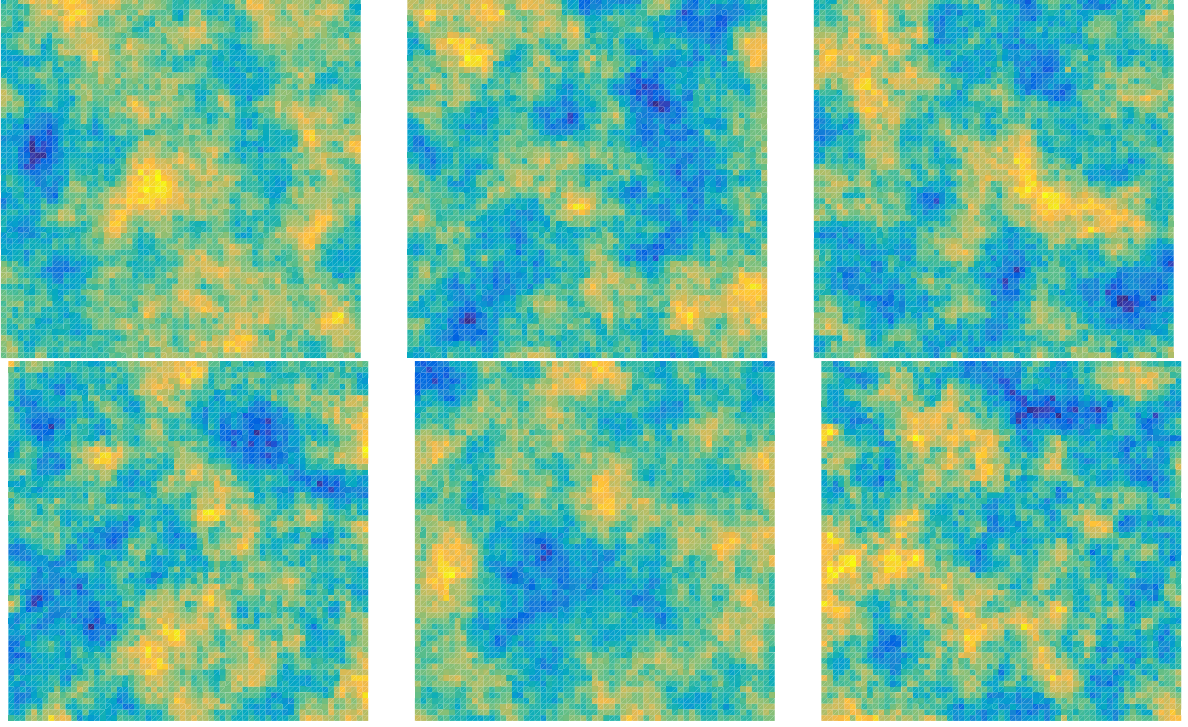


Fig. 2: Samples of \mathcal{Z} with a stationary covariance function from (2) with $p = 2$ and $\mu = 1/2$.

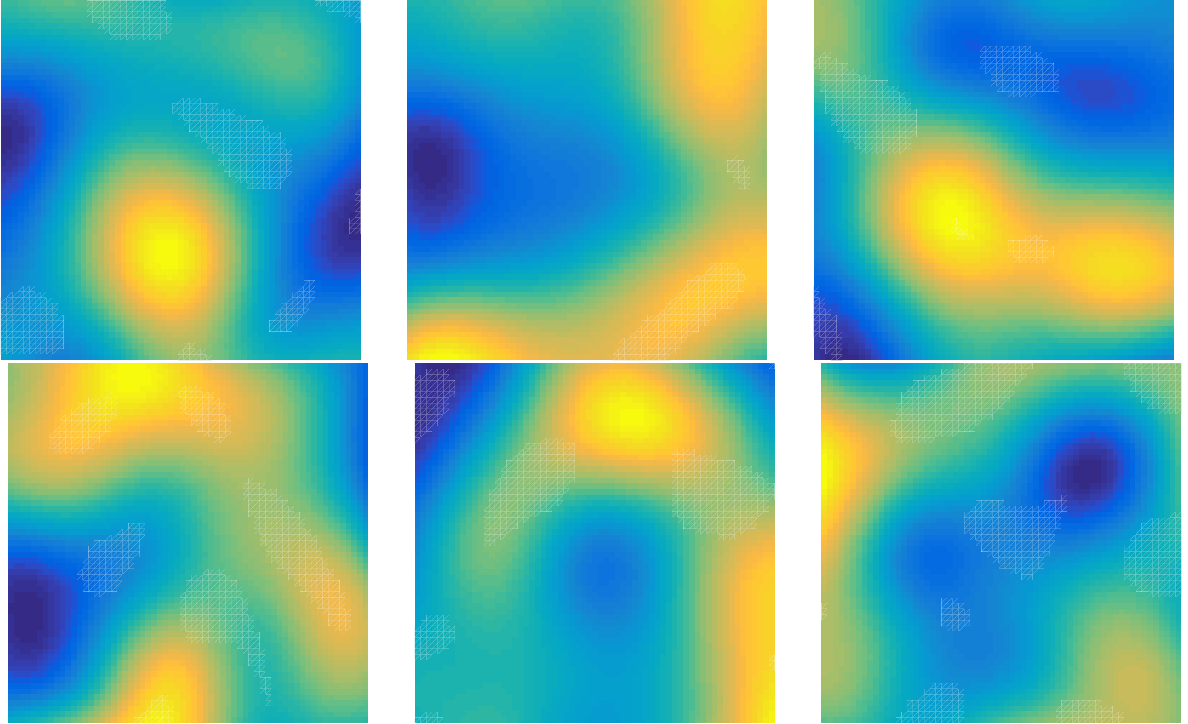


Fig. 3: Samples of \mathcal{Z} with a stationary covariance function from (2) with $p = 2$ and $\mu = \infty$.

computational time, Algorithm 1 is superior in almost all cases (particularly for smooth fields). However, keep in mind that according to Theorem 1, Algorithm 1 needs up to $\mathcal{O}(\log_{\kappa}(\varepsilon)N)$ extra storage, while Algorithm 2 uses only $\mathcal{O}(\log(\log(\varepsilon))N)$ extra storage units (compare Theorem 2).

Figure 5 compares the two algorithms with the direct matrix square root provided by Matlab. We evaluate $\mathcal{Z}(\mathbf{x}, \omega)$ on a Sobol pointset with size 2^m for $m \in \{1, \dots, 14\}$. The number of iterations in both algorithm is set such that the relative error is smaller than 10^{-10} for the example from above with $p = 2$, and varying $\mu \in \{1/2, \infty\}$, $\lambda \in \{1, 10^{-1}, 10^{-2}, 10^{-3}\}$. We see that both, Algorithm 1–2, perform in linear time, whereas the direct approach comes closer to $\mathcal{O}(N^3)$. Even though our H^2 -matrix library is programmed entirely in Matlab (and thus nowhere near optimal performance), the breakthrough point at around $N = 10^3$ shows that also small problems benefit from the speed up.

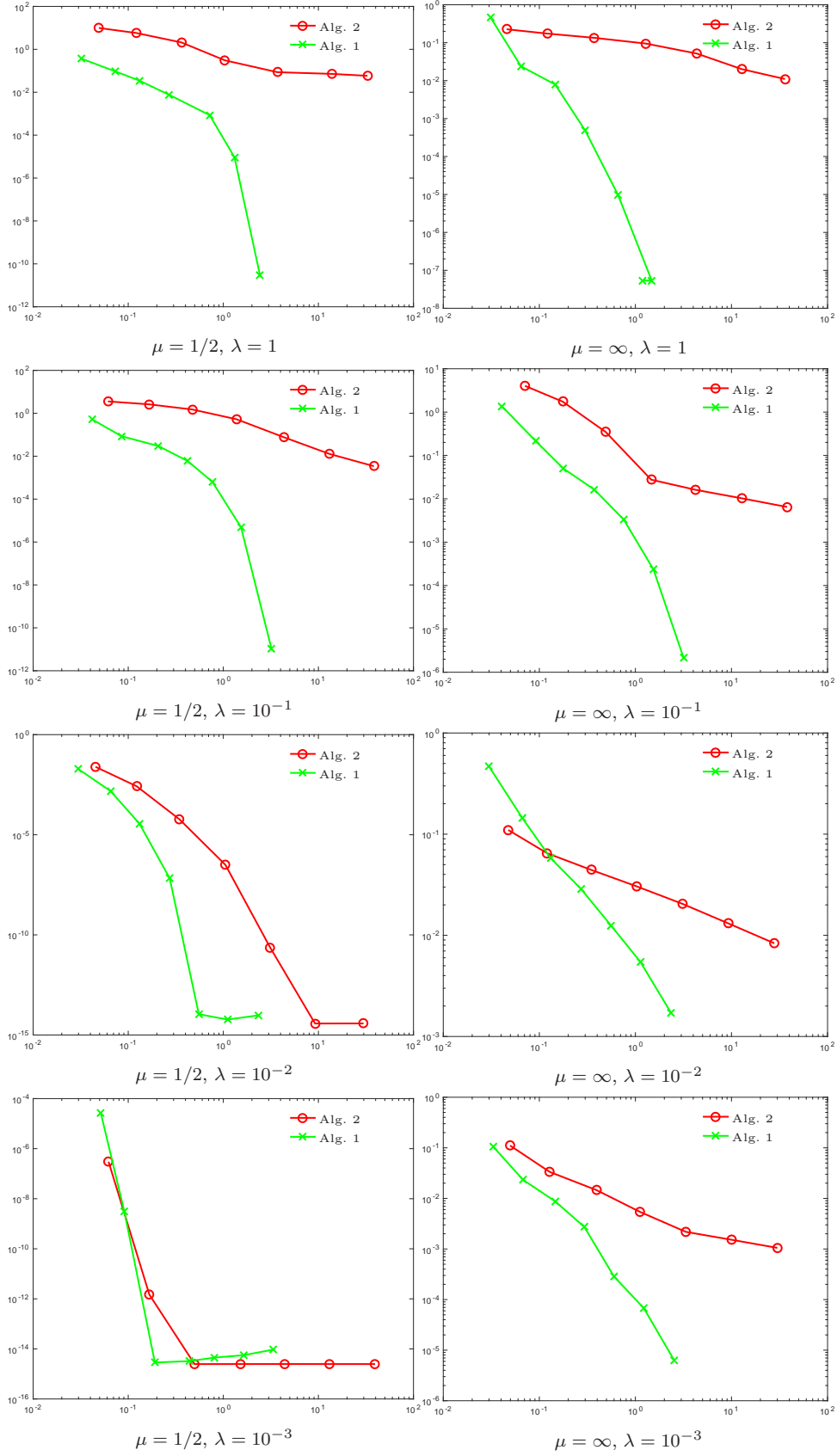


Fig. 4: Comparison of Algorithm 2 and Algorithm 1. We plot the relative error $|\mathcal{Z}_{k,p}(z) - C^{1/2}z|/|z|$ versus computation time in seconds.

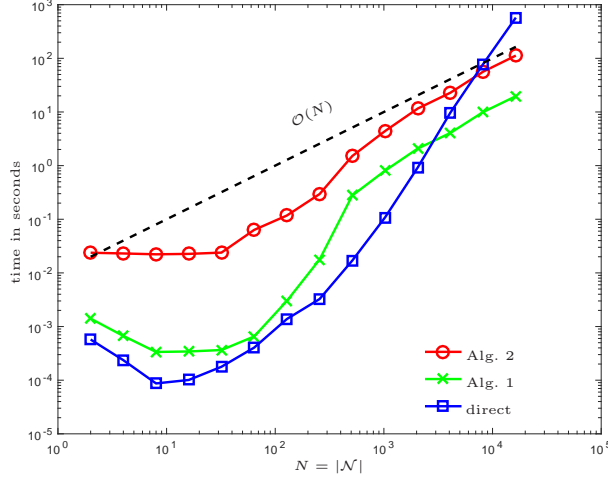


Fig. 5: Computational time in seconds versus the number of evaluation points N . The direct approach uses Matlab's `sqrtrm` function.

5 Lemmas for the proof of Theorem 1

First, we state a well-known result for the convenience of the reader.

Lemma 2 ([17, Lemma 2.2]) *Let $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{N \times N}$ denote positive definite, symmetric matrices. Then, there holds*

$$\|\mathbf{A}^{1/2} - \mathbf{B}^{1/2}\|_2 \leq (\lambda_{\min}(\mathbf{A}) + \lambda_{\min}(\mathbf{B}))^{-1/2} \|\mathbf{A} - \mathbf{B}\|_2.$$

The next lemma quantifies the condition of a truncated Vandermonde matrix. It implies that if the coefficients of the Vandermonde matrix are close to each other in the sense that the right hand side of (16) below is small, then there exists a normalized vector \mathbf{v} such that the matrix vector product with the Vandermonde matrix is small. (Note that the authors suspect that the following result, or parts of it, might have been derived in previous work. However, we could not find any reference and thus include the proof.)

Lemma 3 *Given $M \in \mathbb{N}$, suppose distinct numbers $q_1, \dots, q_M \in \mathbb{R}$ are such that $0 \leq q_i \leq 1$. For $1 \leq k \leq M$, define the transposed Vandermonde-submatrix $\mathbf{V}_k \in \mathbb{R}^{M \times k}$ by*

$$(\mathbf{V}_k)_{ij} := q_i^{j-1}, \quad 1 \leq i \leq M, 1 \leq j \leq k.$$

Then, there exists an absolute constant $C_{\text{van}} > 0$ and $\mathbf{v} \in \mathbb{R}^k$ with $|\mathbf{v}| = 1$ such that $|\mathbf{v}_k| \geq 4^{-k}/(C_{\text{van}}\sqrt{k})$ and

$$|\mathbf{V}_k \mathbf{v}| \leq \sqrt{M} \min_{\substack{SS \subseteq \{1, \dots, M\} \\ |SS|=k}} \max_{1 \leq s \leq M} \prod_{r \in SS \setminus \{t_0(SS)\}} |q_s - q_r|, \quad (16)$$

where $t_0(SS) = \arg \min_{t \in SS} \prod_{r \in SS \setminus \{t\}} |q_t - q_r|$.

Proof Let \mathbf{V} be a $k \times k$ submatrix of \mathbf{V}_k consisting of the rows with labels $SS := \{\nu_1, \dots, \nu_k\} \subseteq \{1, \dots, M\}$ of \mathbf{V}_k are contained in \mathbf{V} . Then, \mathbf{V} is a transposed Vandermonde matrix with coefficients $q_{\nu_1}, \dots, q_{\nu_k}$. Moreover, \mathbf{V} is regular since all the $q_s, s \in SS$ are distinct. It is shown in [16, Equation 2.8.4] (note the transposed definition of Vandermonde matrix there), and indeed is easily verified, that the inverse of \mathbf{V} has the form

$$(\mathbf{V}^{-1})_{ij} = a_{ij}, \quad (17)$$

where the a_{ij} are the coefficients of the Lagrange polynomial

$$L_j(x) := \prod_{\substack{r=1 \\ r \neq j}}^k \frac{x - q_{\nu_r}}{q_{\nu_j} - q_{\nu_r}} = a_{kj}x^{k-1} + a_{(k-1)j}x^{k-2} + \dots + a_{1j}. \quad (18)$$

Expansion of the product in the numerator of L_j shows

$$a_{ij} = \sum_{\substack{\mathcal{R} \subseteq \{1, \dots, k\} \setminus \{j\} \\ |\mathcal{R}| = i-1}} \left(\prod_{r \in \{1, \dots, k\} \setminus (\mathcal{R} \cup \{j\})} (-q_{\nu_r}) \right) \underbrace{\left(\prod_{r \in \{1, \dots, k\} \setminus \{j\}} \frac{1}{q_{\nu_j} - q_{\nu_r}} \right)}_{:= P_{j,k}}. \quad (19)$$

Since there are $\binom{k-1}{i-1}$ subsets of $\{1, \dots, k\} \setminus \{j\}$ with $i-1$ elements, we have

$$|a_{ij}| \leq \binom{k-1}{i-1} \max_{s \in SS} |q_s|^{k-i} |P_{j,k}| \leq \binom{k-1}{\lceil (k-1)/2 \rceil} |P_{j,k}|,$$

where we used $\binom{k}{j} \leq \binom{k}{\lceil k/2 \rceil}$ for all $k, j \in \mathbb{N}$ and $q_s \leq 1$ for all $s \in SS$. With Stirling's formula, we prove that there exists $C > 0$ such that $\binom{k-1}{\lceil (k-1)/2 \rceil} \leq C4^k$ for all $k \in \mathbb{N}$ and thus

$$|a_{ij}| \leq C4^k |P_{j,k}|.$$

Moreover, (19) shows

$$|a_{kj}| = |P_{j,k}|.$$

The last two estimates together with (17) imply for the j -th unit vector \mathbf{e}^j that

$$|\mathbf{V}^{-1} \mathbf{e}^j| = |(a_{1j}, \dots, a_{kj})| \begin{cases} \leq C\sqrt{k}4^k |P_{j,k}|, \\ \geq |a_{kj}| = |P_{j,k}|. \end{cases} \quad (20)$$

Analogously, we get

$$|(\mathbf{V}^{-1} \mathbf{e}^j)_k| = |a_{kj}| = |P_{j,k}|. \quad (21)$$

Define $\mathbf{v}^j := \mathbf{V}^{-1} \mathbf{e}^j / |\mathbf{V}^{-1} \mathbf{e}^j| \in \mathbb{R}^k$, so that $|\mathbf{v}^j| = 1$. From (20)–(21), we see that

$$|(\mathbf{v}^j)_k| \geq 4^{-k} / (C\sqrt{k}),$$

so that we may choose $C_{\text{van}} = C$, as well as

$$|\mathbf{V} \mathbf{v}^j| \leq |P_{j,k}|.$$

Moreover, given any other $k \times k$ submatrix $\tilde{\mathbf{V}}$ of \mathbf{V}_k comprising the rows with labels $\tilde{SS} := \{\mu_1, \dots, \mu_k\} \subseteq \{1, \dots, M\}$ of \mathbf{V}_k , we have by (17)–(18)

$$(\tilde{\mathbf{V}} \mathbf{V}^{-1})_{ij} = \sum_{r=1}^k \tilde{\mathbf{V}}_{ir} (\mathbf{V}^{-1})_{rj} = L_j(q_{\mu_i}) = \prod_{\substack{r=1 \\ r \neq j}}^k \frac{q_{\mu_i} - q_{\nu_r}}{q_{\nu_j} - q_{\nu_r}}.$$

This implies

$$\begin{aligned} |\tilde{\mathbf{V}} \mathbf{v}^j| &= |\tilde{\mathbf{V}} \mathbf{V}^{-1} \mathbf{V} \mathbf{v}^j| \leq \sqrt{k} \max_{1 \leq s, t \leq k} \prod_{\substack{r=1 \\ r \neq t}}^k \frac{|q_{\mu_s} - q_{\nu_r}|}{|q_{\nu_t} - q_{\nu_r}|} |\mathbf{V} \mathbf{v}^j| \\ &\leq \sqrt{k} \max_{1 \leq s, t \leq k} \prod_{\substack{r=1 \\ r \neq t}}^k \frac{|q_{\mu_s} - q_{\nu_r}|}{|q_{\nu_t} - q_{\nu_r}|} \prod_{\substack{r=1 \\ r \neq j}}^k |q_{\nu_j} - q_{\nu_r}|. \end{aligned}$$

We now choose $j = t_0(SS)$ (where the definition of $t_0(\cdot)$ is as in the statement of the lemma). (Note that $t_0(SS)$ is independent of the choice of \tilde{SS}). Then on observing that the maximum over t is attained by $t = t_0(SS)$, we obtain

$$|\tilde{\mathbf{V}} \mathbf{v}^{t_0(SS)}| \leq \sqrt{k} \max_{1 \leq s \leq k} \prod_{\substack{r=1 \\ r \neq t_0(SS)}}^k |q_{\mu_s} - q_{\nu_r}|.$$

Since \widetilde{SS} was an arbitrary subset of cardinality k , we may split the \mathbb{R}^M -vector $\mathbf{V}_k \mathbf{v}^{t_0}$ into a sum of not more than M/k vectors of length k , each of which uses in effect a different subset SS . This shows

$$|\mathbf{V}_k \mathbf{v}^{t_0(SS)}| \leq \sqrt{\frac{M}{k}} \sqrt{k} \max_{1 \leq s \leq M} \prod_{\substack{r=1 \\ r \neq t_0(SS)}}^k |q_s - q_{\nu_r}|.$$

Since SS was an arbitrary subset of cardinality k , we obtain immediately

$$\min_{\substack{SS \subseteq \{1, \dots, M\} \\ |SS|=k}} |\mathbf{V}_k \mathbf{v}^{t_0(SS)}| \leq \sqrt{M} \min_{\substack{SS \subseteq \{1, \dots, M\} \\ |SS|=k}} \max_{1 \leq s \leq M} \prod_{\substack{r=1 \\ r \neq t_0(SS)}}^k |q_s - q_{\nu_r}|$$

and conclude the proof.

The following lemma proves decay properties of the QR -factorization of the Krylov subspace matrix \mathbf{Z} defined below.

Lemma 4 *Let $\mathbf{M} \in \mathbb{R}^{N \times N}$ be symmetric positive definite and assume that $0 < \kappa < 1$ and $C_\kappa > 0$ are such that the sequence of all distinct eigenvalues $\lambda_1 \geq \dots \geq \lambda_M > 0 \in \mathbb{R}$ (for some $M \leq N$) of \mathbf{M} satisfies $|\lambda_i - \lambda_j| \leq \lambda_1 C_\kappa \kappa^{\min\{i,j\}}$ for all $1 \leq i, j \leq M$. Given $1 \leq k \leq M$ and $\mathbf{z} \in \mathbb{R}^N$, define $\mathbf{Z} \in \mathbb{R}^{N \times k}$ by*

$$\mathbf{Z} := (\mathbf{z}, \lambda_1^{-1} \mathbf{M} \mathbf{z}, \lambda_1^{-2} \mathbf{M}^2 \mathbf{z}, \dots, \lambda_1^{-(k-1)} \mathbf{M}^{k-1} \mathbf{z}).$$

Consider the QR -factorization $\mathbf{Z} = \mathbf{Q} \mathbf{R}$, with $\mathbf{Q} \in \mathbb{R}^{N \times k}$ satisfying $\mathbf{Q}^T \mathbf{Q} = \mathbf{I}_N$ and $\mathbf{R} \in \mathbb{R}^{k \times k}$ upper triangular. Then the diagonal entries of \mathbf{R} satisfy

$$|\mathbf{R}_{nn}| \leq C_R \lambda_1^{-1} |\mathbf{z}| \sqrt{Nn} (4C_\kappa)^{n-1} \kappa^{(n-1)n/2} \quad \text{for all } 1 \leq n \leq k. \quad (22)$$

The constant $C_R > 0$ depends only on C_{van} from Lemma 3.

Proof First, we prove the case $n = k$. Let $\mathbf{v}^1, \dots, \mathbf{v}^N \in \mathbb{R}^N$ denote an orthonormal basis of eigenvectors of \mathbf{M} with decreasing eigenvalues $\tilde{\lambda}_1 \geq \dots \geq \tilde{\lambda}_N$ (including multiple eigenvalues). With $\mathbf{z} = \sum_{r=1}^N \mathbf{z}_r \mathbf{v}^r$, there holds

$$\mathbf{M}^j \mathbf{z} = \sum_{r=1}^N \mathbf{z}_r \tilde{\lambda}_r^j \mathbf{v}^r. \quad (23)$$

Define $q_j := \lambda_j / \lambda_1$ for all $j = 1, \dots, M$. Recall the matrix \mathbf{V}_k from Lemma 3 and define the matrix of eigenvectors $\mathbf{B} := (\mathbf{v}^1, \dots, \mathbf{v}^N) \in \mathbb{R}^{N \times N}$. Finally, for $i = 1, \dots, N$, we define $j(i)$ to be the unique $j(i) \in \{1, \dots, M\}$ such that $\lambda_{j(i)} = \tilde{\lambda}_i$. Then, for given $\mathbf{z} \in \mathbb{R}^N$, define the matrix $\mathbf{S} \in \mathbb{R}^{N \times M}$ such that $\mathbf{S}_{ij(i)} := \mathbf{z}_i$ for all $1 \leq i \leq N$ and $\mathbf{S}_{ij} = 0$ if $j \neq j(i)$. From this, because \mathbf{S} contains each component of \mathbf{z} exactly once, we obtain

$$\|\mathbf{S}\|_2 \leq \|\mathbf{S}\|_F = |\mathbf{z}|. \quad (24)$$

As an intermediate step, we want to prove

$$\mathbf{Z} = \mathbf{B} \mathbf{S} \mathbf{V}_k. \quad (25)$$

To that end, note that by definition, we have

$$(\mathbf{B} \mathbf{S})_{ij} = \sum_{\substack{r=1 \\ \tilde{\lambda}_r = \lambda_j}}^N \mathbf{z}_r (\mathbf{v}^r)_i.$$

Hence, multiplication by \mathbf{V}_k shows that

$$\begin{aligned} (\mathbf{BSV}_k)_{ij} &= \sum_{s=1}^M (\mathbf{BS})_{is} (\mathbf{V}_k)_{sj} = \sum_{s=1}^M \sum_{\substack{r=1 \\ \tilde{\lambda}_r = \lambda_s}}^N \mathbf{z}_r(\mathbf{v}^r)_i q_s^{j-1} \\ &= \sum_{s=1}^M \sum_{\substack{r=1 \\ \tilde{\lambda}_r = \lambda_s}}^N \mathbf{z}_r(\mathbf{v}^r)_i \frac{\lambda_s^{j-1}}{\lambda_1^{j-1}} = \lambda_1^{-(j-1)} \sum_{r=1}^N \mathbf{z}_r(\mathbf{v}^r)_i \tilde{\lambda}_r^{j-1} \end{aligned}$$

This and (23) allow us to prove for all $1 \leq i \leq N$ and all $1 \leq j \leq k$

$$\mathbf{Z}_{ij} = (\lambda_1^{-(j-1)} \mathbf{M}^{j-1} \mathbf{z})_i = \sum_{r=1}^N \mathbf{z}_r \frac{\tilde{\lambda}_r^{j-1}}{\lambda_1^{j-1}} \mathbf{v}_i^r = (\mathbf{BSV}_k)_{ij},$$

which proves (25). Thus, and by use of (24), we have with \mathbf{v} from Lemma 3

$$|\mathbf{Z}\mathbf{v}| = |\mathbf{BSV}_k \mathbf{v}| = |\mathbf{SV}_k \mathbf{v}| \leq |\mathbf{z}| \sqrt{M} \min_{\substack{SS \subseteq \{1, \dots, M\} \\ |SS|=k}} \max_{1 \leq s \leq M} \prod_{r \in SS \setminus \{t_0(SS)\}} |q_s - q_r|, \quad (26)$$

where $t_0(SS) = \arg \min_{t \in SS} \prod_{r \in SS \setminus \{t\}} |q_t - q_r|$. By choosing $SS = \{1, \dots, k\}$, we have

$$\min_{\substack{SS \subseteq \{1, \dots, M\} \\ |SS|=k}} \max_{1 \leq s \leq M} \prod_{r \in SS \setminus \{t_0(SS)\}} |q_s - q_r| \leq \max_{1 \leq s \leq M} \prod_{\substack{r=1 \\ r \neq t_0(\{1, \dots, k\})}}^k |q_s - q_r| = \prod_{\substack{r=1 \\ r \neq t_0(SS)}}^k |q_M - q_r|, \quad (27)$$

since for $r \leq k$ each factor $|q_s - q_r|$ is maximized by taking $s = M$. By the decay assumption on the λ_i it follows that

$$\prod_{\substack{r=1 \\ r \neq t_0(\{1, \dots, k\})}}^k |q_M - q_r| \leq \frac{1}{\lambda_1} \prod_{r=1}^{k-1} C_\kappa \kappa^r = \frac{1}{\lambda_1} C_\kappa^{k-1} \kappa^{k(k-1)/2}. \quad (28)$$

The combination of (26)–(28) proves

$$|\mathbf{Z}\mathbf{v}| \leq \frac{1}{\lambda_1} |\mathbf{z}| \sqrt{N} C_\kappa^{k-1} \kappa^{k(k-1)/2}$$

For the QR -factorization, this implies

$$|\mathbf{R}_{kk}| \leq |\mathbf{R}\mathbf{v}|/|\mathbf{v}_k| = |\mathbf{Z}\mathbf{v}|/|\mathbf{v}_k| \leq \frac{1}{\lambda_1} C |\mathbf{z}| \sqrt{Nk} (4C_\kappa)^k \kappa^{k(k-1)/2},$$

where we used the lower bound on \mathbf{v}_k from Lemma 3, and $Q^T Q = I_N$. This concludes the proof for $n = k$. For all $1 \leq n < k$, we may repeat the proof with $\mathbf{Z}|_{\{1, \dots, N\} \times \{1, \dots, n\}}$ and use the fact that the first n columns of \mathbf{R} remain unchanged. This concludes the proof.

The next lemma shows that the matrices \mathbf{Q}_j from Algorithm 1 are strongly tied to the matrices $\mathbf{Z} = \mathbf{Q}\mathbf{R}$ defined in Lemma 4.

Lemma 5 *Given $\mathbf{z} \in \mathbb{R}^N$ and let $\mathbf{M} \in \mathbb{R}^{N \times N}$ be symmetric positive definite. Call Algorithm 1 with \mathbf{M} , \mathbf{z} , and $k \in \mathbb{N}$ to compute $k_0 \leq k$ and \mathbf{R}_j , \mathbf{Q}_j for all $1 \leq j \leq k_0$. Define $\mathbf{Z}, \mathbf{Q}, \mathbf{R}$ satisfying $\mathbf{Z} = \mathbf{Q}\mathbf{R}$ as in Lemma 4. Then, \mathbf{Z} has full rank if and only if $k_0 = k$. In this case, \mathbf{Q}_j (as defined in Algorithm 1) for $1 \leq j \leq k$ satisfies $\mathbf{Q}_j = \mathbf{Q}|_{\{1, \dots, N\} \times \{1, \dots, j\}}$, i.e., the first j columns coincide and*

$$\text{range}(\mathbf{Q}_j) = \text{span}\{\mathbf{z}, \dots, \mathbf{M}^{j-1} \mathbf{z}\} = \text{range}(\mathbf{Z}|_{\{1, \dots, N\} \times \{1, \dots, j\}}) \quad (29)$$

for all $1 \leq j \leq k$.

Proof Let \mathbf{q}^j denote the j -th column of \mathbf{Q}_j . First, assume $k_0 = k$ and hence

$$(\mathbf{R}_j)_{jj} \neq 0 \quad \text{for all } 1 \leq j \leq k. \quad (30)$$

We show

$$\text{range}(\mathbf{Q}_j) = \text{span}\{\mathbf{z}, \dots, \mathbf{M}^{j-1}\mathbf{z}\} \quad (31)$$

for all $1 \leq j \leq k$ by induction. To that end, note that $\mathbf{Q}_1 = \mathbf{q}^1 = \mathbf{z}/|\mathbf{z}|$. Thus, (31) holds for $j = 1$. Assume (31) holds for all $1 \leq j < j_0 \leq k$. By construction of the matrices in Algorithm 1, we have

$$(\mathbf{Q}_{j_0-1}, \mathbf{M}\mathbf{q}^{j_0-1}) = \mathbf{Q}_{j_0}\mathbf{R}_{j_0}. \quad (32)$$

By the induction assumption, $\mathbf{q}^{j_0-1} \in \text{span}\{\mathbf{z}, \dots, \mathbf{M}^{j_0-2}\mathbf{z}\}$. Thus, (32) and the fact that \mathbf{R}_{j_0} is regular (by (30)) imply

$$\text{range}(\mathbf{Q}_{j_0}) = \text{span}\{\text{range}(\mathbf{Q}_{j_0-1}), \mathbf{M}\mathbf{q}^{j_0-1}\} \subseteq \text{span}\{\mathbf{z}, \dots, \mathbf{M}^{j_0-1}\mathbf{z}\}.$$

The fact that \mathbf{Q}_{j_0} is orthogonal (and hence its range is j_0 dimensional) shows even equality, that is

$$\text{range}(\mathbf{Q}_{j_0}) = \text{span}\{\mathbf{z}, \dots, \mathbf{M}^{j_0-1}\mathbf{z}\}. \quad (33)$$

This concludes the induction, and proves (31) for all $1 \leq j \leq k$. The second equation in (29) follows by definition of \mathbf{Z} . Particularly, we see that \mathbf{Z} has full rank, since \mathbf{Q}_k is orthogonal.

From (32), we see by use of the fact that \mathbf{Q}_{j-1} is already orthogonal, that we have $\mathbf{Q}_j|_{\{1, \dots, N\} \times \{1, \dots, j-1\}} = \mathbf{Q}_{j-1}$ for all $2 \leq j \leq k$. Together with (31), this means

$$\text{range}(\mathbf{q}^1, \dots, \mathbf{q}^j) = \text{range}(\mathbf{Q}_j) = \text{span}\{\mathbf{z}, \dots, \mathbf{M}^{j-1}\mathbf{z}\}$$

for all $1 \leq j \leq k$. Since \mathbf{Z} has full rank and \mathbf{R} is upper triangular, we also have

$$\text{range}(\mathbf{Q}|_{\{1, \dots, N\} \times \{1, \dots, j\}}) = \text{span}\{\mathbf{z}, \dots, \mathbf{M}^{j-1}\mathbf{z}\}$$

for all $1 \leq j \leq k$. Since both matrices are orthogonal, there holds

$$\mathbf{Q}|_{\{1, \dots, N\} \times \{1, \dots, j\}} = (\mathbf{q}^1, \dots, \mathbf{q}^j) = \mathbf{Q}_j \quad \text{for all } 1 \leq j \leq k.$$

For the converse implication, assume that \mathbf{Z} has full rank and $k_0 < k$. Therefore, there exists a minimal j_0 such that $(\mathbf{R}_{j_0})_{j_0 j_0} = 0$ (note that $j_0 > 1$ since $\mathbf{R}_{1,11} = 1$ by construction). Then, since $(\mathbf{R}_j)_{jj} \neq 0$ for all $1 \leq j < j_0$, the identity (29) shows $\text{range}(\mathbf{Z}|_{\{1, \dots, N\} \times \{1, \dots, j\}}) = \text{range}(\mathbf{Q}_j)$ for all $j < j_0$. From this, we argue that

$$\mathbf{q}^{j_0-1} \in \text{range}(\mathbf{Z}|_{\{1, \dots, N\} \times \{1, \dots, j_0-1\}}) \setminus \text{range}(\mathbf{Z}|_{\{1, \dots, N\} \times \{1, \dots, j_0-2\}}),$$

which shows $\text{range}((\mathbf{Q}_{j_0-1}, \mathbf{M}\mathbf{q}^{j_0-1})) = \text{range}(\mathbf{Z}|_{\{1, \dots, N\} \times \{1, \dots, j_0\}})$ and therefore the matrix $(\mathbf{Q}_{j_0-1}, \mathbf{M}\mathbf{q}^{j_0-1})$ has full rank. Hence, (32) implies that \mathbf{R}_{j_0} has full rank, which in particular implies $(\mathbf{R}_{j_0})_{j_0 j_0} \neq 0$. This contradicts the assumption on j_0 and thus proves $(\mathbf{R}_j)_{jj} \neq 0$ for all $1 \leq j \leq k_0$.

The following result proves that if Algorithm 1 terminates in less than k steps (due to the criterion in step (1c)), the quantity $\mathbf{M}^{1/2}\mathbf{z}$ is computed exactly.

Lemma 6 *Let $\mathbf{z} \in \mathbb{R}^N$ and let $\mathbf{M} \in \mathbb{R}^{N \times N}$ be symmetric positive definite. Call Algorithm 1 with \mathbf{M} , \mathbf{z} , and $k \in \mathbb{N}$ to compute $k_0 \leq k$ as well as \mathbf{Q}_j for all $1 \leq j \leq k_0$. Define $\mathbf{U}_{k_0} = \mathbf{Q}_{k_0}^T \mathbf{M} \mathbf{Q}_{k_0}$ as in Algorithm 1. If $k_0 < k$, there holds*

$$\mathbf{M}^{1/2}\mathbf{z} = \mathbf{Q}_{k_0} \mathbf{U}_{k_0}^{1/2} \mathbf{Q}_{k_0}^T \mathbf{z}.$$

Proof If $k < k_0$ then Lemma 5 shows that \mathbf{Z} as defined in Lemma 4 does not have full rank. Moreover, Lemma 5 (with k in the call to Algorithm 1 replaced by k_0) shows that $\mathbf{Z}|_{\{1,\dots,N\} \times \{1,\dots,k_0\}}$ has full rank. By definition of \mathbf{Z} , this implies $\text{range}(\mathbf{Z}|_{\{1,\dots,N\} \times \{1,\dots,k_0\}}) = \text{range}(\mathbf{Z})$. Therefore, (29) shows

$$\begin{aligned} \text{range}(\mathbf{M}\mathbf{Q}_{k_0}) &\subseteq \text{range}(\mathbf{M}\mathbf{Z}|_{\{1,\dots,N\} \times \{1,\dots,k_0\}}) \subseteq \text{range}(\mathbf{Z}) \\ &= \text{range}(\mathbf{Z}|_{\{1,\dots,N\} \times \{1,\dots,k_0\}}) = \text{range}(\mathbf{Q}_{k_0}). \end{aligned} \quad (34)$$

Let $\overline{\mathbf{Q}} \in \mathbb{R}^{N \times N}$ be an orthonormal matrix such that its first k_0 columns coincide with \mathbf{Q}_{k_0} , i.e., $\overline{\mathbf{Q}} = (\mathbf{Q}_{k_0}, \mathbf{Q}_\perp)$ for some orthonormal $\mathbf{Q}_\perp \in \mathbb{R}^{N \times (N-k_0)}$. We obtain

$$\mathbf{M}^{1/2} = \overline{\mathbf{Q}} \overline{\mathbf{Q}}^T \mathbf{M}^{1/2} \overline{\mathbf{Q}} \overline{\mathbf{Q}}^T = \overline{\mathbf{Q}} (\overline{\mathbf{Q}}^T \mathbf{M} \overline{\mathbf{Q}})^{1/2} \overline{\mathbf{Q}}^T. \quad (35)$$

There holds

$$\overline{\mathbf{Q}}^T \mathbf{M} \overline{\mathbf{Q}} = \begin{pmatrix} \mathbf{Q}_{k_0}^T \mathbf{M} \mathbf{Q}_{k_0} & \mathbf{Q}_{k_0}^T \mathbf{M} \mathbf{Q}_\perp \\ \mathbf{Q}_\perp^T \mathbf{M} \mathbf{Q}_{k_0} & \mathbf{Q}_\perp^T \mathbf{M} \mathbf{Q}_\perp \end{pmatrix}.$$

The invariance property (34) shows $\mathbf{Q}_\perp^T \mathbf{M} \mathbf{Q}_{k_0} = \mathbf{0}$, and by symmetry also $\mathbf{Q}_{k_0}^T \mathbf{M} \mathbf{Q}_\perp = \mathbf{0}$. Therefore, we have

$$(\overline{\mathbf{Q}}^T \mathbf{M} \overline{\mathbf{Q}})^{1/2} = \begin{pmatrix} \mathbf{U}_{k_0}^{1/2} & \mathbf{0} \\ \mathbf{0} & (\mathbf{Q}_\perp^T \mathbf{M} \mathbf{Q}_\perp)^{1/2} \end{pmatrix}.$$

This and (35), together with $\mathbf{z} \in \text{range}(\mathbf{Q}_{k_0})$, show $\mathbf{M}^{1/2} \mathbf{z} = \overline{\mathbf{Q}} (\overline{\mathbf{Q}}^T \mathbf{M} \overline{\mathbf{Q}})^{1/2} \overline{\mathbf{Q}}^T \mathbf{z} = \mathbf{Q}_{k_0} \mathbf{U}_{k_0}^{1/2} \mathbf{Q}_{k_0}^T \mathbf{z}$ and conclude the proof.

The following result is the main tool to prove Theorem 1 (i).

Lemma 7 *Let $\mathbf{z} \in \mathbb{R}^N$ and let $\mathbf{M} \in \mathbb{R}^{N \times N}$ be symmetric positive definite. Call Algorithm 1 with \mathbf{M} , \mathbf{z} , and $k \in \mathbb{N}$ to compute $k_0 \leq k$ as well as \mathbf{Q}_j for all $1 \leq j \leq k_0$. Let $\mathbf{U}_{k_0} = \mathbf{Q}_{k_0}^T \mathbf{M} \mathbf{Q}_{k_0}$ be defined as in Algorithm 1. Then, there holds*

$$\begin{aligned} &\frac{|\mathbf{M}^{1/2} \mathbf{z} - \mathbf{Q}_{k_0} \mathbf{U}_{k_0}^{1/2} \mathbf{Q}_{k_0}^T \mathbf{z}|}{|\mathbf{z}|} \\ &\leq \begin{cases} \sqrt{2} \sqrt{\lambda_{\min}(\mathbf{M}) + \lambda_{\max}(\mathbf{M})} \left(1 - \frac{2\lambda_{\min}(\mathbf{M})}{\lambda_{\min}(\mathbf{M}) + \lambda_{\max}(\mathbf{M})}\right)^{(k-1) \log_3(2)/2} & \text{if } k_0 = k, \\ 0 & \text{if } k_0 < k. \end{cases} \end{aligned}$$

Proof The case $k_0 < k$ is covered in Lemma 6. Assume $k_0 = k$. Note that $\mathbf{Q}_k \mathbf{Q}_k^T$ is the identity on $\text{range}(\mathbf{Q}_k)$. Lemma 5 shows that $\mathbf{M}^j \mathbf{z} \in \text{range}(\mathbf{Q}_k)$ for all $0 \leq j \leq k-1$. Hence, we have

$$\mathbf{M}^j \mathbf{z} = \mathbf{Q}_k \mathbf{Q}_k^T \mathbf{M}^j \mathbf{Q}_k \mathbf{Q}_k^T \mathbf{z} = \mathbf{Q}_k (\mathbf{Q}_k^T \mathbf{M} \mathbf{Q}_k)^j \mathbf{Q}_k^T \mathbf{z} \quad \text{for all } 1 \leq j \leq k-1.$$

Thus, any polynomial $p \in \mathcal{P}^{k-1}$ of degree $k-1$ satisfies

$$p(\mathbf{M}) \mathbf{z} = \mathbf{Q}_k \mathbf{Q}_k^T p(\mathbf{M}) \mathbf{Q}_k \mathbf{Q}_k^T \mathbf{z} = \mathbf{Q}_k p(\mathbf{Q}_k^T \mathbf{M} \mathbf{Q}_k) \mathbf{Q}_k^T \mathbf{z} = \mathbf{Q}_k p(\mathbf{U}_k) \mathbf{Q}_k^T \mathbf{z}.$$

This implies for all $p \in \mathcal{P}^{k-1}$

$$\begin{aligned} &|\mathbf{M}^{1/2} \mathbf{z} - \mathbf{Q}_k \mathbf{U}_k^{1/2} \mathbf{Q}_k^T \mathbf{z}| \\ &\leq |\mathbf{M}^{1/2} \mathbf{z} - \mathbf{Q}_k p(\mathbf{U}_k) \mathbf{Q}_k^T \mathbf{z}| + |\mathbf{Q}_k p(\mathbf{U}_k) \mathbf{Q}_k^T \mathbf{z} - \mathbf{Q}_k \mathbf{U}_k^{1/2} \mathbf{Q}_k^T \mathbf{z}| \\ &\leq |\mathbf{M}^{1/2} \mathbf{z} - p(\mathbf{M}) \mathbf{z}| + |\mathbf{Q}_k (p(\mathbf{U}_k) - \mathbf{U}_k^{1/2}) \mathbf{Q}_k^T \mathbf{z}| \\ &\leq (\|\mathbf{M}^{1/2} - p(\mathbf{M})\|_2 + \|p(\mathbf{U}_k) - \mathbf{U}_k^{1/2}\|_2) |\mathbf{z}|. \end{aligned} \quad (36)$$

For the first term on the right-hand side, Lemma 9 below proves

$$\|\mathbf{M}^{1/2} - s^{-1/2} \mathbf{A}_j\|_2 \leq s^{-1/2} \left(1 - \frac{2\lambda_{\min}(\mathbf{M})}{\lambda_{\min}(\mathbf{M}) + \lambda_{\max}(\mathbf{M})}\right)^{2^j}$$

for all $j \in \mathbb{N}$, $s = 2/(\lambda_{\min}(\mathbf{M}) + \lambda_{\max}(\mathbf{M}))$ and \mathbf{A}_j is defined in (14). By expanding the recursive definition of the Schulz iteration in (14), we see $s^{-1/2}\mathbf{A}_j = p(\mathbf{M})$ for some $p \in \mathcal{P}^{3^j}$. Since the maximal $j \in \mathbb{N}$ such that $3^j \leq k-1$ satisfies $2^j \geq (k-1)^{\log_3(2)}/2$, this proves that

$$\|\mathbf{M}^{1/2} - p(\mathbf{M})\|_2 \leq s^{-1/2} \left(1 - \frac{2\lambda_{\min}(\mathbf{M})}{\lambda_{\min}(\mathbf{M}) + \lambda_{\max}(\mathbf{M})}\right)^{(k-1)^{\log_3(2)}/2}.$$

Since \mathbf{U}_k is an orthogonal projection of \mathbf{M} , we have $\lambda_{\min}(\mathbf{M}) \leq \lambda_{\min}(\mathbf{U}_k) \leq \lambda_{\max}(\mathbf{U}_k) \leq \lambda_{\max}(\mathbf{M})$. Thus, repeating the above argument for \mathbf{U}_k instead of \mathbf{M} yields

$$\|\mathbf{M}^{1/2} - p(\mathbf{M})\|_2 + \|p(\mathbf{U}_k) - \mathbf{U}_k^{1/2}\|_2 \leq 2s^{-1/2} \left(1 - \frac{2\lambda_{\min}(\mathbf{M})}{\lambda_{\min}(\mathbf{M}) + \lambda_{\max}(\mathbf{M})}\right)^{(k-1)^{\log_3(2)}/2}.$$

This in combination with (36) and Lemma 5 conclude the proof.

This result shows that the subspace $\text{range}(\mathbf{Q}_j)$ is almost invariant under \mathbf{M} .

Lemma 8 *Let $\mathbf{z} \in \mathbb{R}^N$ and let $\mathbf{M} \in \mathbb{R}^{N \times N}$ be symmetric positive definite. Call Algorithm 1 with \mathbf{M} , \mathbf{z} , and $k \in \mathbb{N}$ to compute $k_0 \leq k$ as well as \mathbf{Q}_j for all $1 \leq j \leq k_0$. Let \mathbf{q}^j be the last column of \mathbf{Q}_j for all $1 \leq j \leq k_0$. There holds for all $1 \leq j < k_0$*

$$\|\mathbf{M}\mathbf{Q}_j - \mathbf{Q}_j\mathbf{Q}_j^T\mathbf{M}\mathbf{Q}_j\|_2 = |(\mathbf{q}^{j+1})^T\mathbf{M}\mathbf{q}^j| \quad (37)$$

as well as

$$\min_{1 \leq i \leq j} \|\mathbf{M}\mathbf{Q}_i - \mathbf{Q}_i\mathbf{Q}_i^T\mathbf{M}\mathbf{Q}_i\|_2 \leq 8|\lambda_1|C_R^{1/j}N^{1/(2j)}C_\kappa\kappa^{(j+1)/2}.$$

Proof Recall $\mathbf{Z}, \mathbf{Q}, \mathbf{R}$ satisfying $\mathbf{Z} = \mathbf{Q}\mathbf{R}$ from Lemma 4 with k replaced by k_0 in the call to Algorithm 1. By Lemma 5, \mathbf{q}^j coincides with the j -th column of \mathbf{Q} for all $1 \leq j \leq k_0$. Moreover, let \mathbf{r}^j be the j -th column of \mathbf{R} and let $\tilde{\mathbf{r}}^j$ be the j -th column of \mathbf{R}^{-1} from Lemma 4. All quantities are well-defined since \mathbf{Z} has maximal rank k_0 by Lemma 5. For the first statement (37), note that $\text{range}(\mathbf{M}\mathbf{Q}_j) \subseteq \text{range}(\mathbf{Q}_{j+1})$ implies $\mathbf{M}\mathbf{Q}_j = \mathbf{Q}_{j+1}\mathbf{Q}_{j+1}^T\mathbf{M}\mathbf{Q}_j$. Moreover, due to Lemma 5, we have $\mathbf{Q}_{j+1} = (\mathbf{Q}_j, \mathbf{q}^{j+1})$, and hence $\mathbf{Q}_{j+1}\mathbf{Q}_{j+1}^T = \mathbf{q}^{j+1}(\mathbf{q}^{j+1})^T + \mathbf{Q}_j\mathbf{Q}_j^T$. Altogether, this shows

$$\begin{aligned} \|\mathbf{M}\mathbf{Q}_j - \mathbf{Q}_j\mathbf{Q}_j^T\mathbf{M}\mathbf{Q}_j\|_2 &= \|\mathbf{Q}_{j+1}\mathbf{Q}_{j+1}^T\mathbf{M}\mathbf{Q}_j - \mathbf{Q}_j\mathbf{Q}_j^T\mathbf{M}\mathbf{Q}_j\|_2 \\ &= \|(\mathbf{q}^{j+1}(\mathbf{q}^{j+1})^T + \mathbf{Q}_j\mathbf{Q}_j^T)\mathbf{M}\mathbf{Q}_j - \mathbf{Q}_j\mathbf{Q}_j^T\mathbf{M}\mathbf{Q}_j\|_2 \\ &= \|\mathbf{q}^{j+1}(\mathbf{q}^{j+1})^T\mathbf{M}\mathbf{Q}_j\|_2 = |(\mathbf{q}^{j+1})^T\mathbf{M}\mathbf{q}^j|, \end{aligned}$$

where the last step follows because \mathbf{q}^{j+1} is orthogonal to $\mathbf{M}\mathbf{q}^i$, $i = 1, \dots, j-1$. This proves (37).

To see the remaining statement, note that the definition of \mathbf{Z} implies

$$(\mathbf{M}\mathbf{Z})|_{\{1, \dots, N\} \times \{j\}} = \lambda_1 \mathbf{Z}|_{\{1, \dots, N\} \times \{j+1\}} = \lambda_1 (\mathbf{Q}\mathbf{R})|_{\{1, \dots, N\} \times \{j+1\}} = \lambda_1 \mathbf{Q}\mathbf{r}^{j+1}$$

as well as

$$\mathbf{q}^j = (\mathbf{Z}\mathbf{R}^{-1})|_{\{1, \dots, N\} \times \{j\}} = \mathbf{Z}\tilde{\mathbf{r}}^j.$$

The last two identities, and the fact (see above) that $(\mathbf{q}^{j+1})^T\mathbf{M}\mathbf{Z}|_{\{1, \dots, N\} \times \{i\}} = 0$ for all $1 \leq i \leq j-1$, imply

$$\begin{aligned} (\mathbf{q}^{j+1})^T\mathbf{M}\mathbf{q}^j &= (\mathbf{q}^{j+1})^T\mathbf{M}\mathbf{Z}\tilde{\mathbf{r}}^j = (\mathbf{q}^{j+1})^T\mathbf{M}\mathbf{Z}|_{\{1, \dots, N\} \times \{j\}}(\tilde{\mathbf{r}}^j)_j \\ &= \lambda_1(\mathbf{q}^{j+1})^T\mathbf{Q}\mathbf{r}^{j+1}(\tilde{\mathbf{r}}^j)_j = \lambda_1(\mathbf{r}^{j+1})_{j+1}(\tilde{\mathbf{r}}^j)_j. \end{aligned}$$

The triangular structure of \mathbf{R} implies $(\mathbf{R}^{-1})_{jj} = 1/\mathbf{R}_{jj}$ and hence $(\tilde{\mathbf{r}}^j)_j = 1/\mathbf{R}_{jj}$ (where $\mathbf{R}_{jj} \neq 0$ by assumption). This shows

$$(\mathbf{q}^{j+1})^T\mathbf{M}\mathbf{q}^j = \lambda_1 \frac{\mathbf{R}_{(j+1)(j+1)}}{\mathbf{R}_{jj}}. \quad (38)$$

With Lemma 4, we have

$$\begin{aligned} \frac{|\mathbf{R}_{(j+1)(j+1)}|}{|\mathbf{R}_{jj}|} \frac{|\mathbf{R}_{jj}|}{|\mathbf{R}_{(j-1)(j-1)}|} \cdots \frac{|\mathbf{R}_{22}|}{|\mathbf{R}_{11}|} |\mathbf{R}_{11}| &= |\mathbf{R}_{(j+1)(j+1)}| \\ &\leq C_R |\mathbf{z}| \sqrt{N} \sqrt{j+1} (4C_\kappa)^j \kappa^{j(j+1)/2}. \end{aligned} \quad (39)$$

Moreover, we know $|\mathbf{R}_{11}| = |\mathbf{q}^1 \mathbf{R}_{11}| = |\mathbf{z}|$. This implies that at least one of the fractions on the left-hand side of (39) must be smaller than the j -th root of the right hand side of (39) divided by $|\mathbf{z}|$ and hence

$$\min_{1 \leq i \leq j} \frac{|\mathbf{R}_{(i+1)(i+1)}|}{|\mathbf{R}_{ii}|} \leq C_R^{1/j} N^{1/(2j)} (j+1)^{1/j} 4C_\kappa \kappa^{(j+1)/2}.$$

With this, (38), and (37), we obtain

$$\min_{1 \leq i \leq j} \|\mathbf{M}\mathbf{Q}_i - \mathbf{Q}_i \mathbf{Q}_i^T \mathbf{M}\mathbf{Q}_i\|_2 \leq 8|\lambda_1| C_R^{1/j} N^{1/(2j)} C_\kappa \kappa^{(j+1)/2},$$

where we used $\sqrt{j+1}^{1/j} \leq 2$. This concludes the proof.

The following proposition is the main tool to prove Theorem 1 (ii).

Proposition 2 *Let $\mathbf{z} \in \mathbb{R}^N$ and let $\mathbf{M} \in \mathbb{R}^{N \times N}$ be symmetric positive definite. Call Algorithm 1 with \mathbf{M} , \mathbf{z} , and $k \in \mathbb{N}$ to compute $k_0 \leq k$ as well as \mathbf{Q}_j for all $1 \leq j \leq k_0$. Then, $\mathbf{U}_j := \mathbf{Q}_j^T \mathbf{M}\mathbf{Q}_j$ satisfies the error bound*

$$\frac{|(\mathbf{M}^{1/2})\mathbf{z} - \mathbf{Q}_j(\mathbf{U}_j^{1/2})\mathbf{Q}_j^T \mathbf{z}|}{|\mathbf{z}|} \leq \begin{cases} \frac{|(\mathbf{q}^{j+1})^T \mathbf{M}\mathbf{q}^j|}{\sqrt{\lambda_{\min}(\mathbf{M})}} & 1 \leq j < k_0, \\ 0 & j = k_0 \text{ and } k_0 < k \end{cases}$$

and we have the a priori estimate

$$\min_{1 \leq i \leq j} |(\mathbf{q}^{i+1})^T \mathbf{M}\mathbf{q}^i| \leq 8\lambda_1 C_R^{1/j} N^{1/(2j)} C_\kappa \kappa^{(j+1)/2}$$

for all $1 \leq j < k_0$.

Proof The case $k_0 < k$ and $j = k_0$ is trivially covered in Lemma 6. For the other cases, let $\overline{\mathbf{Q}} \in \mathbb{R}^{N \times N}$ be orthonormal such that the first j columns coincide with \mathbf{Q}_j , i.e., $\overline{\mathbf{Q}} = (\mathbf{Q}_j, \mathbf{Q}_\perp)$ for some orthonormal $\mathbf{Q}_\perp \in \mathbb{R}^{N \times (N-j)}$. Then, we write

$$\overline{\mathbf{Q}}^T \mathbf{M} \overline{\mathbf{Q}} = \begin{pmatrix} \mathbf{U}_j & \mathbf{R} \\ \mathbf{0} & \mathbf{T} \end{pmatrix}$$

for matrices $\mathbf{R} = \mathbf{Q}_\perp^T \mathbf{M}\mathbf{Q}_j \in \mathbb{R}^{j \times (N-j)}$, $\mathbf{T} \in \mathbb{R}^{(N-j) \times (N-j)}$. This means that

$$\left\| \overline{\mathbf{Q}}^T \mathbf{M} \overline{\mathbf{Q}} - \begin{pmatrix} \mathbf{U}_j & \mathbf{0} \\ \mathbf{0} & \mathbf{T} \end{pmatrix} \right\|_2 \leq \|\mathbf{R}\|_2 =: R_j.$$

Lemma 2 then implies

$$\left\| (\overline{\mathbf{Q}}^T \mathbf{M} \overline{\mathbf{Q}})^{1/2} - \begin{pmatrix} \mathbf{U}_j^{1/2} & \mathbf{0} \\ \mathbf{0} & \mathbf{T}^{1/2} \end{pmatrix} \right\|_2 \leq \lambda_{\min}(\mathbf{M})^{-1/2} R_j \quad (40)$$

Since $\mathbf{I} - \mathbf{Q}_j \mathbf{Q}_j^T = \mathbf{Q}_\perp \mathbf{Q}_\perp^T$, we have

$$R_j = \|\mathbf{R}\|_2 = \|\mathbf{Q}_\perp^T \mathbf{M}\mathbf{Q}_j\|_2 = \|\mathbf{Q}_\perp \mathbf{Q}_\perp^T \mathbf{M}\mathbf{Q}_j\|_2 = \|\mathbf{M}\mathbf{Q}_j - \mathbf{Q}_j \mathbf{Q}_j^T \mathbf{M}\mathbf{Q}_j\|_2.$$

With $(\bar{Q}^T M \bar{Q})^{1/2} = \bar{Q}^T M^{1/2} \bar{Q}$ and since the ranges of Q_j and Q_\perp are orthogonal, we have $Q_\perp^T Q_j Q_j^T = 0$ and

$$\begin{aligned}
& \|M^{1/2} Q_j Q_j^T - Q_j (U_j^{1/2}) Q_j^T\|_2 \\
&= \|M^{1/2} Q_j Q_j^T - Q_j (U_j^{1/2}) Q_j^T Q_j Q_j^T - Q_\perp (T^{1/2}) Q_\perp^T Q_j Q_j^T\|_2 \\
&\leq \|M^{1/2} - Q_j (U_j^{1/2}) Q_j^T - Q_\perp (T^{1/2}) Q_\perp^T\|_2 \\
&= \|\bar{Q}^T (M^{1/2} - Q_j (U_j^{1/2}) Q_j^T - Q_\perp (T^{1/2}) Q_\perp^T) \bar{Q}\|_2 \\
&= \left\| (\bar{Q}^T M \bar{Q})^{1/2} - \begin{pmatrix} U_j^{1/2} & 0 \\ 0 & T^{1/2} \end{pmatrix} \right\|_2.
\end{aligned} \tag{41}$$

The combination of (40) and (41) shows

$$\|M^{1/2} Q_j Q_j^T - Q_j (U_j^{1/2}) Q_j^T\|_2 \leq \lambda_{\min}(M)^{-1/2} R_j.$$

We conclude the proof with $z = Q_j Q_j^T z$ due to $z \in \text{range}(Q_j)$ and Lemma 8.

6 Lemma for the proof of Theorem 2

The following lemma is the main tool for the proof of Theorem 2.

Lemma 9 *Let $M \in \mathbb{R}^{N \times N}$ be symmetric and positive definite. Then, the iteration (14) with initial values $A_0 = sM$ and $B_0 = I$ satisfies*

$$\|M^{1/2} - s^{-1/2} A_k\|_2 \leq s^{-1/2} (\max\{|1 - s\lambda_{\max}(M)|, |1 - s\lambda_{\min}(M)|\})^{2^k} \tag{42}$$

for all $k \in \mathbb{N}$ and all $s > 0$. The minimum bound is attained at $s = 2/(\lambda_{\min}(M) + \lambda_{\max}(M))$ such that $\max\{|1 - s\lambda_{\max}(M)|, |1 - s\lambda_{\min}(M)|\} = 1 - 2\lambda_{\min}(M)/(\lambda_{\min}(M) + \lambda_{\max}(M))$.

Proof Straightforward calculations show

$$\left(\|M^{1/2} - A_k\|_2^2 + \|M^{-1/2} - B_k\|_2^2 \right)^{1/2} = \left\| \begin{pmatrix} 0 & M^{1/2} \\ M^{-1/2} & 0 \end{pmatrix} - \begin{pmatrix} 0 & A_k \\ B_k & 0 \end{pmatrix} \right\|_2.$$

We want to employ [14, Theorem 5.2]. In their notation, setting $k = 1$ (which is different to our k) and $m = 0$ in [14, Theorem 5.2], we show

$$\left\| \begin{pmatrix} 0 & M^{1/2} \\ M^{-1/2} & 0 \end{pmatrix} - \begin{pmatrix} 0 & A_k \\ B_k & 0 \end{pmatrix} \right\|_2 < \left\| \begin{pmatrix} I - M & 0 \\ 0 & I - M \end{pmatrix} \right\|_2^{2^k} = \|I - M\|_2^{2^k}$$

for all $k \in \mathbb{N}$. By scaling of M , we may minimize the right-hand side. To that end, we observe that the spectrum satisfies $\sigma(I - sM) \subset [1 - s\lambda_{\max}(M), 1 - s\lambda_{\min}(M)]$. The fact $\|I - sM\|_2 \leq \max\{|1 - s\lambda_{\max}(M)|, |1 - s\lambda_{\min}(M)|\}$ proves (42). A straightforward optimization of $s > 0$ concludes the proof.

A Proof of Lemma 1

The following lemma is an elementary statement on holomorphic functions

Lemma 10 *Let $f: O \rightarrow \mathbb{C}$ be a continuous function on the domain $O \subset \mathbb{C}^n$ which is holomorphic in O in all variables x_i , $i \in \{1, \dots, n\}$, i.e.,*

$$x_i \mapsto f(x_1, \dots, x_i, \dots, x_n)$$

is holomorphic in $\{x_i \in \mathbb{C} : (x_1, \dots, x_i, \dots, x_n) \in O\}$ for all $x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n \in \mathbb{C}$. Then, for all multi-indices $\alpha \in \mathbb{N}_0^n$, the function $\partial_x^\alpha f$ is holomorphic in O in all variables x_i , $i \in \{1, \dots, n\}$ as defined above.

Proof The result is proved by induction on $|\alpha|_1$. Obviously, for $|\alpha|_1 = 0$, $\partial_{\mathbf{x}}^{\alpha} f = f$ and the statement is true. Assume the statement holds for all $|\alpha|_1 \leq k$ and choose some $\alpha \in \mathbb{N}_0^n$ with $|\alpha|_1 = k + 1$. Then, we have for some $i \in \{1, \dots, n\}$ and some $\alpha_0 \in \mathbb{N}_0^n$ with $|\alpha_0|_1 = k$ that

$$\partial_{\mathbf{x}}^{\alpha} f = \partial_{\mathbf{x}_i} \partial_{\mathbf{x}}^{\alpha_0} f.$$

Since, $\partial_{\mathbf{x}}^{\alpha_0} f$ is holomorphic in O in all variables by the induction hypothesis, obviously $\partial_{\mathbf{x}}^{\alpha} f$ is holomorphic in O at least in \mathbf{x}_i (derivatives of holomorphic functions are holomorphic). To prove the statement for all other variables, we may employ Cauchy's integral formula to obtain

$$\partial_{\mathbf{x}}^{\alpha} f(\mathbf{x}) = \partial_{\mathbf{x}_i} \partial_{\mathbf{x}}^{\alpha_0} f = \frac{1}{2\pi i} \int_{\partial B_{\varepsilon}(\mathbf{x}_i)} \frac{\partial_{\mathbf{x}}^{\alpha_0} f(\mathbf{x}_1, \dots, \mathbf{x}_{i-1}, \mathbf{z}, \mathbf{x}_{i+1}, \dots, \mathbf{x}_n)}{(\mathbf{z} - \mathbf{x}_i)^2} d\mathbf{z},$$

for some $\varepsilon > 0$ with $B_{\varepsilon}(\mathbf{x}_i) \subset \mathbb{C}$ being the ball with radius ε . The integrand is holomorphic in all variables \mathbf{x}_j , $j \neq i$. Hence, we conclude that $\partial_{\mathbf{x}}^{\alpha} f(\mathbf{x})$ is holomorphic in all variables and prove the assertion.

The following results is elementary but technical.

Lemma 11 For $n, p \in \mathbb{N}$, define the set $M := \{\mathbf{x} \in \mathbb{C}^n : \text{real}(\sum_{i=1}^n \mathbf{x}_i^p) \leq 0\}$. Then, there holds $(\mathbb{R}^n)_+ := \{\mathbf{x} \in \mathbb{R}^n \setminus \{0\} : \mathbf{x}_i \geq 0\} \cap M = \emptyset$ and

$$\text{dist}(M, \mathbf{x}) \geq |\sin(\frac{\pi}{2p})| |\mathbf{x}| \quad \text{for all } \mathbf{x} \in (\mathbb{R}^n)_+.$$

Proof Let $\mathbf{x} \in (\mathbb{R}^n)_+$, then we have $\sum_{i=1}^n \mathbf{x}_i^p > 0$ and hence $\mathbf{x} \notin M$. It is easy to see that the cone $C_p := \{r \exp(i\phi) : r > 0, \phi \in (-\frac{\pi}{2p}, \frac{\pi}{2p})\} \subset \mathbb{C}$ satisfies $\text{real}(x^p) > 0$ for all $x \in C_p$. Thus, we have that

$$C_p^n := \left(\prod_{i=1}^n (\{0\} \cup C_p) \right) \setminus \{0\} \subset \mathbb{C}^n$$

satisfies $C_p^n \cap M = \emptyset$. Moreover, a simple geometric argument shows that all $x > 0$ satisfy

$$\text{dist}(x, \partial C_p) = x \sin(\pi/(2p)).$$

Since $(\mathbb{R}^n)_+ \subseteq C_p^n$, this implies

$$\text{dist}(M, \mathbf{x}) \geq \text{dist}(\partial C_p^n, \mathbf{x}) = \left(\sum_{i=1}^n \mathbf{x}_i^2 \sin(\pi/(2p))^2 \right)^{1/2} = |\sin(\pi/(2p))| |\mathbf{x}|.$$

This concludes the proof.

Products of asymptotically smooth functions are again asymptotically smooth. This is shown in the next lemma.

Lemma 12 Given two functions $f, g: D \times D \rightarrow \mathbb{R}$ which are asymptotically smooth (1). Then, also their product fg satisfies (1).

Proof To simplify the notation, we consider f, g as functions of one variable $\mathbf{z} = (\mathbf{x}, \mathbf{y}) \in D \times D \subset \mathbb{R}^{2d}$. For multi-indices $\alpha, \beta \in \mathbb{N}^{2d}$, define

$$\binom{\alpha}{\beta} := \prod_{i=1}^{2d} \binom{\alpha_i}{\beta_i}.$$

Note that there holds $\binom{\alpha}{\beta} \leq \binom{|\alpha|_1}{|\beta|_1}$. This follows from the basic combinatorial fact that the number of possible choices of β_i elements out of a set of α_i elements for all $i = 1, \dots, 2d$ is smaller than the number of choices of $|\beta|_1$ elements out of a set of $|\alpha|_1$ elements.

The Leibniz formula together with the definition of asymptotically smooth function (1) show for $\alpha \in \mathbb{N}^{2d}$

$$\begin{aligned} |\partial_{\mathbf{z}}^{\alpha} fg(\mathbf{z})| &\leq \sum_{\substack{\beta \in \mathbb{N}_0^{2d} \\ \beta \leq \alpha}} \binom{\alpha}{\beta} |\partial_{\mathbf{z}}^{\beta} f(\mathbf{z})| |\partial_{\mathbf{z}}^{\alpha-\beta} g(\mathbf{z})| \\ &\leq \sum_{\substack{\beta \in \mathbb{N}_0^{2d} \\ \beta \leq \alpha}} \binom{|\alpha|_1}{|\beta|_1} c_1(c_2|\mathbf{x} - \mathbf{y}|)^{-|\beta|_1} |\beta|_1! c_1(c_2|\mathbf{x} - \mathbf{y}|)^{|\alpha|_1 - |\beta|_1} (|\alpha|_1 - |\beta|_1)! \\ &\leq \sum_{\substack{\beta \in \mathbb{N}_0^{2d} \\ \beta \leq \alpha}} c_1^2(c_2|\mathbf{x} - \mathbf{y}|)^{-|\alpha|_1} |\alpha|_1! \leq (|\alpha|_1 + 1)^{2d} c_1^2(c_2|\mathbf{x} - \mathbf{y}|)^{-|\alpha|_1} |\alpha|_1! \\ &\lesssim c_1^2(c_2|\mathbf{x} - \mathbf{y}|)^{-|\alpha|_1} |\alpha|_1!. \end{aligned}$$

This concludes the proof.

The final lemma of this section proves the concatenations of certain asymptotically smooth functions are asymptotically smooth.

Lemma 13 *Let $g: D \times D \rightarrow \mathbb{R}$ be asymptotically smooth (1) with constants $c_1, c_2 > 0$.*

- (i) *If $c_g := \sup_{\mathbf{x} \in D \times D} g(\mathbf{x}) < \infty$. Then, $\exp \circ g$ satisfies (1) with constants $\tilde{c}_1 := \exp(c_g)$ and $\tilde{c}_2 := c_2/(2 \max\{1, c_1\})$.*
- (ii) *If g satisfies $\partial_{\mathbf{x}}^\alpha \partial_{\mathbf{y}}^\beta g(\mathbf{x}, \mathbf{y}) \leq C_g$ for all $\alpha, \beta \in \mathbb{N}_0^d$ and some $C_g < \infty$ as well as $g(\mathbf{x}, \mathbf{y}) \geq C_g^{-1} |\mathbf{x} - \mathbf{y}|$, then, $g^{1/q}$ satisfies (1) with $\tilde{\varrho}_1 = 1/2$ and $\tilde{\varrho}_2 = C_g^{-1}$ for all $q \in \mathbb{N}$.*
- (iii) *If g satisfies the assumptions from (ii) and additionally $g(\mathbf{x}, \mathbf{y}) \geq c_0 > 0$ for all $\mathbf{x}, \mathbf{y} \in D$, then $g^{-1/q}$ satisfies (1) for all $q \in \mathbb{N}$.*

Proof To simplify the notation, we consider g as a function of one variable $\mathbf{z} = (\mathbf{x}, \mathbf{y}) \in D \times D \subset \mathbb{R}^{2d}$. Define the set of all partitions of $\{1, \dots, n\}$ as

$$\Pi(n) := \left\{ P \subseteq 2^{\{1, \dots, n\}} : S \cap S' = \emptyset, S, S' \in P, \bigcup_{S \in P} S = \{1, \dots, n\} \right\}.$$

For (i), Faà di Bruno's formula shows for all multi indices $\alpha \in \mathbb{N}^{2d}$ with $n = |\alpha|_1$ that

$$\partial_{\mathbf{z}}^\alpha (\exp \circ g)(\mathbf{z}) = \sum_{P \in \Pi(n)} \exp \circ g(\mathbf{z}) \prod_{S \in P} \left(\prod_{j \in S} \partial_{\mathbf{z}_{\alpha_j}} \right) g(\mathbf{z}),$$

where $\partial_{\mathbf{z}_{\alpha_j}} g(\mathbf{z}) := 0$ for all $j > 2d$. The definition of asymptotically smooth (1) implies

$$\begin{aligned} |\partial_{\mathbf{z}}^\alpha (\exp \circ g)(\mathbf{z})| &\leq \exp(c_g) \sum_{P \in \Pi(n)} \prod_{S \in P} c_1 (c_2 |\mathbf{x} - \mathbf{y}|)^{-|S|} |S|! \\ &\leq \exp(c_g) \max\{1, c_1\}^n (c_2 |\mathbf{x} - \mathbf{y}|)^{-n} \sum_{P \in \Pi(n)} \prod_{S \in P} |S|!. \end{aligned}$$

With $f(x) := (1 - x)^{-1}$, $x \in \mathbb{R}$, the last factor can be written, using Faà di Bruno's formula again, as

$$\sum_{P \in \Pi(n)} \prod_{S \in P} |S|! = \sum_{P \in \Pi(n)} \exp \circ f(0) \prod_{S \in P} \partial_{\mathbf{x}}^{|S|} f(0) = \partial_x^n (\exp((1 - x)^{-1}))(0).$$

As the function $h(x) := \exp((1 - x)^{-1})$, $x \in \mathbb{C}$ is holomorphic at least for $|x| < 1$, Cauchy's integral formula shows

$$|\partial_x^n h(0)| = \frac{n!}{2\pi} \left| \int_{|z|=1/2} \frac{h(z)}{z^{n+1}} dz \right| \leq n! 2^n \exp(2).$$

Altogether, we conclude the proof of (i) by

$$|\partial_{\mathbf{z}}^\alpha (\exp \circ g)(\mathbf{z})| \leq \exp(c_g) \left(\frac{c_2}{2 \max\{1, c_1\}} |\mathbf{x} - \mathbf{y}| \right)^{-n} n!.$$

For (ii), statement, Faà di Bruno's formula shows again

$$|\partial_{\mathbf{z}}^\alpha (g^{1/q})(\mathbf{z})| \leq \sum_{P \in \Pi(n)} |P|! |g(\mathbf{z})|^{1/q - |P|} \prod_{S \in P} C_g \leq C_g^n |\mathbf{x} - \mathbf{y}|^{-|n|} \sum_{P \in \Pi(n)} |P|!.$$

With $r(x) := \exp(x) - 1$, $x \in \mathbb{R}$, the last factor satisfies

$$\sum_{P \in \Pi(n)} |P|! = \sum_{P \in \Pi(n)} (\partial_x^{|P|} f) \circ r(0) \prod_{S \in P} (\partial_x^{|S|} r)(0) = \partial_x^n (f \circ r)(0).$$

The function $h(x) := f \circ r(x) = (2 - \exp(x))^{-1}$, $x \in \mathbb{C}$ is holomorphic at least for $|x| \leq 1/2$. As above, this implies

$$\partial_x^n (f \circ r)(0) \leq n! 2^n$$

and thus concludes the proof of (ii).

For (iii), we conclude the proof as for (ii) by use of the estimate $g(z)^{1/q - |P|} \leq c_0^{-1 - n}$.

At last, we are ready to prove Lemma 1.

Proof (Proof of Lemma 1) To see (1), consider $\varrho(\cdot, \cdot)$ from (2). We define for complex variables $\mathbf{x}_i, \mathbf{y}_i \in \mathbb{C}$

$$d(\mathbf{x} - \mathbf{y}) = \left(\sum_{i=1}^d (\mathbf{x}_i - \mathbf{y}_i)^p \right)^{1/p} \in \mathbb{C}$$

and consider $\tilde{\varrho}(\mathbf{x}, \mathbf{y})$ which is $\varrho(\mathbf{x}, \mathbf{y})$ from (2) but with $d(\mathbf{x} - \mathbf{y})$ instead of $|\mathbf{x} - \mathbf{y}|_p$. With the notation of Lemma 11, the above sum has positive real part in $O := \{(\mathbf{x}, \mathbf{y}) \in \mathbb{C}^{2d} : \mathbf{x} - \mathbf{y} \notin M\}$. Thus, the function $(\mathbf{x}, \mathbf{y}) \mapsto d(\mathbf{x} - \mathbf{y})$ is holomorphic in each variable in O . Since for $a > 0$, $\mathbf{x} \mapsto \mathbf{x}^\mu K_\mu(a\mathbf{x})$ is a holomorphic function on $\mathbb{C} \setminus (\mathbb{R}_- \cup \{0\})$, and $d(\mathbf{x} - \mathbf{y})$ has positive

real part, we deduce that $(\mathbf{x}, \mathbf{y}) \mapsto \tilde{\varrho}(\mathbf{x}, \mathbf{y})$ is holomorphic in each variable in O . Thus, Lemma 10 proves that $\partial_{\mathbf{x}}^{\alpha} \partial_{\mathbf{y}}^{\beta} \tilde{\varrho}(\mathbf{x}, \mathbf{y})$ is holomorphic in O in all variables \mathbf{x}_i and \mathbf{y}_i . Therefore, Cauchy's integral formula applied in all variables shows

$$\begin{aligned} \partial_{\mathbf{x}}^{\alpha} \partial_{\mathbf{y}}^{\beta} \tilde{\varrho}(\mathbf{x}, \mathbf{y}) &= \frac{\prod_{i=1}^d \alpha_i! \beta_i!}{(2\pi i)^{2d}} \int_{\partial B_{\mathbf{x},1}} \cdots \int_{\partial B_{\mathbf{x},d}} \int_{\partial B_{\mathbf{y},1}} \cdots \int_{\partial B_{\mathbf{y},d}} \frac{\tilde{\varrho}(s, t)}{\prod_{i=1}^d (s_i - \mathbf{x}_i)^{\alpha_i+1} (t_i - \mathbf{y}_i)^{\beta_i+1}} dt ds. \end{aligned}$$

The balls $B_{\mathbf{x},i}$ and $B_{\mathbf{y},i}$ have to be chosen such that $\prod_{i=1}^d B_{\mathbf{x},i} \times \prod_{i=1}^d B_{\mathbf{y},i} \subset O$. With Lemma 11, and for $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{2d}$ such that $\mathbf{x} - \mathbf{y} \in (\mathbb{R}^n)_+$ (note that Lemma 11 implies $(\mathbf{x}, \mathbf{y}) \in O$), this can be achieved by setting $B_{\mathbf{x},i} := B_{\varepsilon}(\mathbf{x}_i)$ and $B_{\mathbf{y},i} := B_{\varepsilon}(\mathbf{y}_i)$ with $\varepsilon := \sin(\pi/(2p))|\mathbf{x} - \mathbf{y}|/(2d+1)$. From this, we obtain the estimate

$$|\partial_{\mathbf{x}}^{\alpha} \partial_{\mathbf{y}}^{\beta} \varrho(\mathbf{x}, \mathbf{y})| = |\partial_{\mathbf{x}}^{\alpha} \partial_{\mathbf{y}}^{\beta} \tilde{\varrho}(\mathbf{x}, \mathbf{y})| \lesssim \frac{\alpha! \beta! (2d+1)^{|\alpha|_1 + |\beta|_1}}{|\mathbf{x} - \mathbf{y}|^{|\alpha|_1 + |\beta|_1}} \max_{(s,t) \in D \times D} |\tilde{\varrho}(s, t)| \quad (43)$$

for all $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{2d}$ such that $\mathbf{x} - \mathbf{y} \in (\mathbb{R}^n)_+$, where the first equality follows from $d(\mathbf{x} - \mathbf{y}) = |\mathbf{x} - \mathbf{y}|_p$ for all $\mathbf{x} - \mathbf{y} \in (\mathbb{R}^n)_+$. To remove the restriction $\mathbf{x} - \mathbf{y} \in (\mathbb{R}^n)_+$, consider $b \in \{0, 1\}^d$ and define the function

$$F_b(\mathbf{x}, \mathbf{y}) := ((-1)^{b_1} \mathbf{x}_1, \dots, (-1)^{b_d} \mathbf{x}_d, (-1)^{b_1} \mathbf{y}_1, \dots, (-1)^{b_d} \mathbf{y}_d).$$

Obviously, there holds $\varrho \circ F_b = \varrho$. Since for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ with $\mathbf{x} \neq \mathbf{y}$, there exists some $b \in \{0, 1\}^d$ such that $(\mathbf{x}_b, \mathbf{y}_b) := F_b(\mathbf{x}, \mathbf{y})$ satisfies $\mathbf{x}_b - \mathbf{y}_b \in (\mathbb{R}^n)_+$, we prove (43) for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ with $\mathbf{x} \neq \mathbf{y}$. Finally, the fact $\alpha! \beta! \leq |\alpha + \beta|_1!$, proves that $\varrho(\cdot, \cdot)$ from (2) is asymptotically smooth (1).

Next, consider the covariance function $\varrho(\cdot, \cdot)$ from (3). By definition $\Sigma_{\mathbf{x}}$ is continuous on \overline{D} . Hence, $\det(\Sigma_{\mathbf{x}}) \geq c_0 > 0$ for all $\mathbf{x} \in D$. The assumption (4) implies that also $\det(\Sigma_{\mathbf{x}})$ has bounded derivatives in the sense of (4) (since $\det(\Sigma_{\mathbf{x}})$ is a polynomial in the matrix entries of $\Sigma_{\mathbf{x}}$). Thus, Lemma 13 shows that the functions $(\mathbf{x}, \mathbf{y}) \mapsto \det(\Sigma_{\mathbf{x}})^{1/4}$, $(\mathbf{x}, \mathbf{y}) \mapsto \det(\Sigma_{\mathbf{y}})^{1/4}$, and $(\mathbf{x}, \mathbf{y}) \mapsto \det(\Sigma_{\mathbf{x}} + \Sigma_{\mathbf{y}})^{-q}$, $q \in \{1/2, 1\}$ satisfy (1). With $\Sigma_{\mathbf{x}}$, also all functions $\tilde{\Sigma}_{\mathbf{x}}$ defined by considering only sub-matrices of $\Sigma_{\mathbf{x}}$ satisfy (4). Thus, Cramer's rule and Lemma 12 show that the map $(\mathbf{x}, \mathbf{y}) \mapsto ((\Sigma_{\mathbf{x}} + \Sigma_{\mathbf{y}})^{-1})_{i,j}$ for all $i, j \in \{1, \dots, d\}$ satisfies (1). From this, we conclude (again with Lemma 12), that $(\mathbf{x}, \mathbf{y}) \mapsto (\mathbf{x} - \mathbf{y})^T (\Sigma_{\mathbf{x}} + \Sigma_{\mathbf{y}})^{-1} (\mathbf{x} - \mathbf{y})$ as sum and product of asymptotically smooth functions is asymptotically smooth (1). Finally, Lemma 13 shows that $\varrho(\mathbf{x}, \mathbf{y})$ satisfies (1). This concludes the proof.

B Proof of Proposition 1

The following lemmas state facts about the H^2 -matrix block partitioning, which are well-known but cannot be found explicitly in the literature.

Lemma 14 *Under Assumption 1, there exists a constant $C_B > 0$ which depends only on d and $B_{X_{\text{root}}}$ such that all $X \in \mathbb{T}_{\text{cl}}$ satisfy*

$$\text{diam}(B_X)^d \leq C_B |B_X|, \quad (44a)$$

$$C_B^{-1} (1 + N |B_X| / |D|) \leq |X| \leq C_B (1 + N |B_X| / |D|), \quad (44b)$$

$$|B_X| = 2^{-\text{level}(X)} |B_{X_{\text{root}}}|. \quad (44c)$$

Moreover, all $(X, Y) \in \mathbb{T}$ satisfy

$$C_{BB}^{-1} \text{diam}(B_X) \leq \text{diam}(B_Y) \leq C_{BB} \text{diam}(B_X), \quad (45)$$

where $C_{BB} > 0$ depends only on C_B and C_u .

Proof The first estimate (44a) follows from the fact that always the longest edge of a bounding box is halved. This means that the ratio L_{\max}/L_{\min} of the maximal and the minimal side length of a bounding box B_X stays bounded in terms of the corresponding ratio for $B_{X_{\text{root}}}$. Therefore, we have

$$\text{diam}(B_X)^d \leq (\sqrt{d} L_{\max})^d \lesssim d^{d/2} L_{\min}^d \leq d^{d/2} |B_X|.$$

To see the second estimate (44b), consider a given bounding box B with side lengths L_1, \dots, L_d . It contains less than

$$m_B := \prod_{i=1}^d (L_i + \max_{\mathbf{x} \in \mathcal{N}} \text{diam}(Q_{\mathbf{x}})) / \min_{\mathbf{x} \in \mathcal{N}} |Q_{\mathbf{x}}|$$

points of \mathcal{N} . Due to Assumption 1, there holds

$$\prod_{i=1}^d (L_i + \max_{\mathbf{x} \in \mathcal{N}} \text{diam}(Q_{\mathbf{x}})) \lesssim \prod_{i=1}^d L_i + \prod_{i=1}^d \max_{\mathbf{x} \in \mathcal{N}} \text{diam}(Q_{\mathbf{x}}) \lesssim |B| + |Q_{\mathbf{x}}|,$$

as well as

$$\min_{\mathbf{x} \in \mathcal{N}} |Q_{\mathbf{x}}| \geq C_u^{-1} |D| N^{-1}.$$

Hence, we have the estimate $m_B \lesssim 1 + N|B|/|D|$. Analogously, we derive the converse estimate to prove (44b).

The estimate (44c) follows from the fact $|B_X| = |B_{X'}|/2$ for all $X \in \text{sons}(X')$.

For (45), we observe with (44b)–(44c) that

$$|X| \simeq 2^{-\text{level}(X)}$$

for all $X \in \mathbb{T}_{\text{cl}}$. Thus, for all leaves $X \in \mathbb{T}_{\text{cl}}$ with $\text{sons}(X) = \emptyset$, we have

$$2^{-\text{level}(X)} \simeq C_{\text{leaf}}.$$

By definition of the block-tree \mathbb{T} , a level difference between X and Y for $(X, Y) \in \mathbb{T}$ can only happen, if $\text{sons}(X) = \emptyset$ or $\text{sons}(Y) = \emptyset$. Assume $\text{sons}(X) = \emptyset$. In this case, we have $\text{level}(Y) \geq \text{level}(X)$. Then, we have

$$2^{-\text{level}(X)} \simeq C_{\text{leaf}} \lesssim 2^{-\text{level}(Y)}$$

which implies $\text{level}(X) \geq \text{level}(Y) + C$ for some constant $C > 0$ which depends only on C_{leaf} and C_B from (44). From this we derive (45) by use of (44).

Lemma 15 *Given the definition of \mathbb{T}_{far} in Section 3.1, there exists a constant $C > 0$ such that all $(X, Y) \in \mathbb{T}_{\text{far}}$ satisfy*

$$C^{-1} \text{diam}(B_X) \leq \eta \text{dist}(B_X, B_Y) \leq C \text{diam}(B_X). \quad (46)$$

Proof By definition of the block-partitioning, for any $(X, Y) \in \mathbb{T}_{\text{far}}$ with X' being the father of X (i.e. $X \in \text{sons}(X')$) there holds that B_X, B_Y satisfy (5) and $B_{X'}, B_Y$ do not satisfy (5). We distinguish two cases: First, let

$$\max\{\text{diam}(B_X), \text{diam}(B_Y)\} \leq \eta/2 \text{dist}(B_X, B_Y). \quad (47)$$

With $\text{dist}(B_{X'}, B_Y) \geq \text{dist}(B_X, B_Y) - \text{diam}(B_X)$, we conclude

$$\max\{\text{diam}(B_X), \text{diam}(B_Y)\} \leq \eta \text{dist}(B_{X'}, B_Y).$$

Since $B_{X'}, B_Y$ do not satisfy (5), this implies

$$\begin{aligned} \text{diam}(B_X) &\leq \max\{\text{diam}(B_X), \text{diam}(B_Y)\} \leq \eta \text{dist}(B_{X'}, B_Y) \\ &< \max\{\text{diam}(B_{X'}), \text{diam}(B_Y)\} = \text{diam}(B_{X'}). \end{aligned}$$

By definition, we have $\text{diam}(B_X) \leq \text{diam}(B_{X'}) \leq 2\text{diam}(B_X)$. This shows

$$\text{diam}(B_X) \leq \eta \text{dist}(B_{X'}, B_Y) < 2 \text{diam}(B_X).$$

Again, using $\text{dist}(B_{X'}, B_Y) \geq \text{dist}(B_X, B_Y) - \text{diam}(B_X)$, we prove (46).

Second, let (47) be false. With (5), this implies

$$\eta/2 \text{dist}(B_X, B_Y) < \max\{\text{diam}(B_X), \text{diam}(B_Y)\} \leq \eta \text{dist}(B_X, B_Y).$$

With (45), we prove (46) immediately.

The following lemma gives some basic facts about tensorial Cebyshev-interpolation.

Lemma 16 *Let $f: B \rightarrow \mathbb{R}$ for an axis parallel box $B \subseteq \mathbb{R}^{2d}$ such that $\partial_j^k f \in L^\infty(B)$ for all $j = 1, \dots, d$ and all $0 \leq k \leq p$. Then, the tensorial Cebyshev-interpolation operator of order p , $I_p: C(B) \rightarrow \mathcal{P}^p(B)$ satisfies*

$$\sup_{\mathbf{x} \in B} |I_p f(\mathbf{x}) - f(\mathbf{x})| \leq 2d\Lambda_p^{2d-1} 4 \frac{4^{-p}}{p!} \text{diam}(B)^p \sum_{i=1}^{2d} \|\partial_{\mathbf{x}_i}^p f\|_{L^\infty(B)}, \quad (48)$$

where

$$\Lambda_p := \sup_{f \in C([-1,1])} \frac{\|I_p^{\mathbf{x}} f\|_{L^\infty([-1,1])}}{\|f\|_{L^\infty([-1,1])}} \leq \frac{2}{\pi} \log(p) + 1 \quad (49)$$

is the operator norm of the one dimensional Cebyshev interpolation operator

Proof It is well-known that the one dimensional Cebaysev interpolation operator $I_p^{\mathbf{x}}$ satisfies the error estimate for any $f \in C([-1, 1])$

$$\|u - I_p^{\mathbf{x}} f\|_{L^\infty([-1, 1])} \leq 4 \frac{2^{-p}}{p!} \|\partial^p f\|_{L^\infty([-1, 1])}$$

with an operator norm given in (49). Consider $B := [-1, 1]^{2d}$. Then, there holds with $I_p^{\mathbf{x}_i}$ denoting interpolation in the \mathbf{x}_i -variable $i \in \{1, \dots, 2d\}$

$$\begin{aligned} |f - I_p f| &= |f - I_p^{\mathbf{x}_1} f + I_p^{\mathbf{x}_1} f - I_p^{\mathbf{x}_2} I_p^{\mathbf{x}_1} f + \dots - I_p f| \\ &\leq \sum_{i=1}^{2d} 4 \frac{2^{-p}}{p!} \|\partial_{\mathbf{x}_i}^p I_p^{\mathbf{x}_1} (I_p^{\mathbf{x}_2} \dots I_p^{\mathbf{x}_{i-1}}) f\|_{L^\infty(B)} \leq \sum_{i=1}^{2d} \Lambda_p^{i-1} 4 \frac{2^{-p}}{p!} \|\partial_{\mathbf{x}_i}^p f\|_{L^\infty(B)} \\ &\leq 2d \Lambda_p^{2d-1} 4 \frac{2^{-p}}{p!} \|\partial_{\mathbf{x}_i}^p f\|_{L^\infty(B)}. \end{aligned}$$

Since, for any affine transformation $A: \mathbb{R}^{2d} \rightarrow \mathbb{R}^{2d}$, we have $I_p(f \circ A) = I_p(f) \circ A$, a standard scaling argument concludes the proof.

Proof (Proof of Proposition 1)

We start by proving that $\lambda_{\min}(\mathbf{C}_p) > 0$ if p satisfies (8). To that end, note

$$\begin{aligned} \lambda_{\min}(\mathbf{C}_p) &= \min_{\mathbf{z} \in \mathbb{R}^N \setminus \{0\}} \frac{(\mathbf{C}_p \mathbf{z})^T \mathbf{z}}{|\mathbf{z}|} \geq \min_{\mathbf{z} \in \mathbb{R}^N \setminus \{0\}} \frac{(\mathbf{C} \mathbf{z})^T \mathbf{z}}{|\mathbf{z}|} - \sup_{\mathbf{z} \in \mathbb{R}^N \setminus \{0\}} \frac{((\mathbf{C}_p - \mathbf{C}) \mathbf{z})^T \mathbf{z}}{|\mathbf{z}|} \\ &\geq \lambda_{\min}(\mathbf{C}) - \|\mathbf{C} - \mathbf{C}_p\|_2 \geq \lambda_{\min}(\mathbf{C}) - \|\mathbf{C} - \mathbf{C}_p\|_F, \end{aligned}$$

since the Frobenius norm is an upper bound for the spectral norm. By use of (7) (which is proved below) and (8), we conclude $\lambda_{\min}(\mathbf{C}_p) > 0$.

To see (7), we first estimate the maximal depth of the tree \mathbb{T}_{cl} . With (44b)–(44c), we obtain $C_{\text{leaf}} \leq |X| \lesssim 2^{-\text{level}(X)}$ for all $X \in \mathbb{T}_{\text{cl}}$ with $\text{sons}(X) \neq \emptyset$. Thus, there holds

$$\max_{X \in \mathbb{T}_{\text{cl}}} \text{level}(X) \lesssim \log(|\mathcal{N}|).$$

Second, we aim to bound the so-called sparsity constant

$$\begin{aligned} C_{\text{sparse}} &:= \max_{X \in \mathbb{T}_{\text{cl}}} \left(|\{Y \in \mathbb{T}_{\text{cl}} : (X, Y) \in \mathbb{T}_{\text{near}} \cup \mathbb{T}_{\text{far}}\}| \right. \\ &\quad \left. + |\{Y \in \mathbb{T}_{\text{cl}} : (Y, X) \in \mathbb{T}_{\text{near}} \cup \mathbb{T}_{\text{far}}\}| \right). \end{aligned}$$

The combination of (45) and (46) (from Lemma 15) shows that $(X, Y) \in \mathbb{T}_{\text{far}}$ only if B_Y touches the (hyper-) annulus with center B_X and radii $C^{-1} \text{diam}(B_X)$ and $C \text{diam}(B_X)$. By comparing the volumes of this annulus and of B_Y and using the fact that all the bounding boxes are disjoint, we see that the number of Y such that $(X, Y) \in \mathbb{T}_{\text{far}}$ is bounded in terms of C and the constants in (44).

For $Y \in \mathbb{T}_{\text{cl}}$ such that $(X, Y) \in \mathbb{T}_{\text{near}}$, we have with (44)–(45)

$$\text{diam}(B_X) \simeq \max\{\text{diam}(B_X), \text{diam}(B_Y)\} > \eta \text{dist}(B_X, B_Y).$$

Again, comparing the volumes of the ball with radius $\text{diam}(B_X)$ and of B_Y , we see that the number of Y such that $(X, Y) \in \mathbb{T}_{\text{near}}$ is bounded in terms of the constants in (44). Altogether, we bound C_{sparse} uniformly in terms of the constants of Lemma 14. Now, [2, Lemma 3.38] proves the estimate for storage requirements and [2, Theorem 3.4.2] proves the estimate for matrix-vector multiplication.

It remains to prove the error estimate. To that end, note that since the near field \mathbb{T}_{near} is stored exactly, there holds

$$\|\mathbf{C} - \mathbf{C}_p\|_F^2 = \sum_{(X, Y) \in \mathbb{T}_{\text{far}}} \|\mathbf{C}|_{I(X) \times I(Y)} - V^X M^{XY} (W^Y)^T\|_F^2.$$

Given, $(i, j) \in I(X) \times I(Y)$, we have with the interpolation operator I_p from Lemma 16 and (1)

$$\begin{aligned} |\mathbf{C}_{ij} - (\mathbf{C}_p)_{ij}| &= |\varrho(\mathbf{x}_i, \mathbf{x}_j) - \sum_{n, m=1}^{p^d} \varrho(q_n^X, q_m^Y) L_n^X(\mathbf{x}_i) L_m^Y(\mathbf{x}_j)| = |\varrho(\mathbf{x}_i, \mathbf{x}_j) - (I_p c)(\mathbf{x}_i, \mathbf{x}_j)| \\ &\lesssim (\log(p) + 1)^{2d-1} \frac{4^{-p}}{p!} \text{diam}(B_X \times B_Y)^p \sum_{i=1}^d (\|\partial_{\mathbf{x}_i}^p c\|_{L^\infty(B)} + \|\partial_{\mathbf{y}_i}^p c\|_{L^\infty(B_X \times B_Y)}) \\ &\lesssim (\log(p) + 1)^{2d-1} \frac{4^{-p}}{p!} \text{diam}(B_X \times B_Y)^p (c_2 \text{dist}(B_X, B_Y))^{-p} p!. \end{aligned}$$

With the admissibility condition (5), we get

$$\text{diam}(B_X \times B_Y) \lesssim \max\{\text{diam}(B_X), \text{diam}(B_Y)\} \leq \eta \text{dist}(B_X, B_Y)$$

and hence

$$|\mathbf{C}_{ij} - (\mathbf{C}_p)_{ij}| \lesssim (\log(p) + 1)^{2d-1} \left(\frac{\eta}{4c_2}\right)^p.$$

The combination of the above estimates concludes the proof.

References

1. I. Babuška, B. Andersson, P. J. Smith, and K. Levin. Damage analysis of fiber composites. I. Statistical analysis on fiber scale. *Comput. Methods Appl. Mech. Engrg.*, 172(1-4):27–77, 1999.
2. Steffen Börm. *Efficient numerical methods for non-local operators*, volume 14 of *EMS Tracts in Mathematics*. European Mathematical Society (EMS), Zürich, 2010.
3. Grace Chan and Andrew T.A. Wood. Algorithm as 312: An algorithm for simulating stationary gaussian random fields. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 46(1):171–181, 1997.
4. C. R. Dietrich and G. N. Newsam. Fast and exact simulation of stationary Gaussian processes through circulant embedding of the covariance matrix. *SIAM J. Sci. Comput.*, 18(4):1088–1107, 1997.
5. J. Dölz, H. Harbrecht, and Ch. Schwab. Covariance regularity and h-matrix approximation for rough random fields. *Numerische Mathematik*, pages 1–27, 2016.
6. I. Elishakoff, editor. *Whys and hows in uncertainty modelling*, volume 388 of *CISM Courses and Lectures*. Springer-Verlag, Vienna, 1999. Probability, fuzziness and anti-optimization.
7. Andreas Frommer. Monotone convergence of the Lanczos approximations to matrix functions of Hermitian matrices. *Electron. Trans. Numer. Anal.*, 35:118–128, 2009.
8. I.G. Graham, F.Y. Kuo, D. Nuyens, R. Scheichl, and I.H. Sloan. Quasi-Monte Carlo methods for elliptic PDEs with random coefficients and applications. *Journal of Computational Physics*, 230(10):3668 – 3694, 2011.
9. Wolfgang Hackbusch. *Hierarchical matrices: algorithms and analysis*, volume 49 of *Springer Series in Computational Mathematics*. Springer, Heidelberg, 2015.
10. Helmut Harbrecht, Michael Peters, and Markus Siebenmorgen. Efficient approximation of random fields for numerical applications. *Numer. Linear Algebra Appl.*, 22(4):596–617, 2015.
11. D. Higdon, J. Swall, and J. Kern. Non-stationary spatial modeling.
12. Nicholas J. Higham. Computing real square roots of a real matrix. *Linear Algebra Appl.*, 88/89:405–430, 1987.
13. Nicholas J. Higham. Stable iterations for the matrix square root. *Numer. Algorithms*, 15(2):227–242, 1997.
14. Charles Kenney and Alan J. Laub. Rational iterative methods for the matrix sign function. *SIAM J. Matrix Anal. Appl.*, 12(2):273–291, 1991.
15. Igor Moret. Rational Lanczos approximations to the matrix square root and related functions. *Numer. Linear Algebra Appl.*, 16(6):431–445, 2009.
16. William H. Press, Saul A. Teukolsky, William T. Vetterling, and Brian P. Flannery. *Numerical recipes*. Cambridge University Press, Cambridge, third edition, 2007. The art of scientific computing.
17. Bernhard A. Schmitt. Perturbation bounds for matrix square roots and pythagorean sums. *Linear Algebra and its Applications*, 174:215 – 227, 1992.